



Entropic Optimal Transport between Unbalanced Gaussian Measures has a Closed Form

Hicham Janati, Boris Muzellec, Gabriel Peyré, Marco Cuturi

► To cite this version:

Hicham Janati, Boris Muzellec, Gabriel Peyré, Marco Cuturi. Entropic Optimal Transport between Unbalanced Gaussian Measures has a Closed Form. NeurIPS 2020 - 34th Conference on Neural Information Processing Systems, Dec 2020, Vancouver / Virtuel, Canada. hal-03063834

HAL Id: hal-03063834

<https://hal.science/hal-03063834>

Submitted on 14 Dec 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Entropic Optimal Transport between Unbalanced Gaussian Measures has a Closed Form

Hicham Janati
Inria Saclay
Paris-Saclay, France
hicham.janati@inria.fr

Boris Muzellec
ENSAE,
Paris-Saclay, France
boris.muzellec@ensae.fr

Gabriel Peyré
CNRS and ENS, PSL University
Paris, France
gabriel.peyre@ens.fr

Marco Cuturi
Google Brain, ENSAE
Paris Saclay, France
cuturi@google.com

Abstract

Although optimal transport (OT) problems admit closed form solutions in a very few notable cases, e.g. in 1D or between Gaussians, these closed forms have proved extremely fecund for practitioners to define tools inspired from the OT geometry. On the other hand, the numerical resolution of OT problems using entropic regularization has given rise to many applications, but because there are no known closed-form solutions for entropic regularized OT problems, these approaches are mostly algorithmic, not informed by elegant closed forms. In this paper, we propose to fill the void at the intersection between these two schools of thought in OT by proving that the entropy-regularized optimal transport problem between two Gaussian measures admits a closed form. Contrary to the unregularized case, for which the explicit form is given by the Wasserstein-Bures distance, the closed form we obtain is differentiable everywhere, even for Gaussians with degenerate covariance matrices. We obtain this closed form solution by solving the fixed-point equation behind Sinkhorn’s algorithm, the default method for computing entropic regularized OT. Remarkably, this approach extends to the generalized *unbalanced* case — where Gaussian measures are scaled by positive constants. This extension leads to a closed form expression for unbalanced Gaussians as well, and highlights the mass transportation / destruction trade-off seen in unbalanced optimal transport. Moreover, in both settings, we show that the optimal transportation plans are (scaled) Gaussians and provide analytical formulas of their parameters. These formulas constitute the first non-trivial closed forms for entropy-regularized optimal transport, thus providing a ground truth for the analysis of entropic OT and Sinkhorn’s algorithm.

1 Introduction

Optimal transport (OT) theory [49, 21] has recently inspired several works in data science, where dealing with and comparing probability distributions, and more generally positive measures, is an important staple (see [41] and references therein). For these applications of OT to be successful, a belief now widely shared in the community is that some form of regularization is needed for OT to be both scalable and avoid the curse of dimensionality [18, 22]. Two approaches have emerged in recent years to achieve these goals: either regularize directly the measures themselves, by looking at them through a simplified lens; or regularize the original OT problem using various modifications. The first approach exploits well-known closed-form identities for OT when comparing two univariate

measures or two multivariate Gaussian measures. In this approach, one exploits those formulas and operates by summarizing complex measures as one or possibly many univariate or multivariate Gaussian measures. The second approach builds on the fact that for arbitrary measures, regularizing the OT problem, either in its primal or dual form, can result in simpler computations and possibly improved sample complexity. The latter approach can offer additional benefits for data science: because the original marginal constraints of the OT problem can also be relaxed, regularized OT can also yield useful tools to compare measures with different total mass — the so-called “unbalanced” case [3]— which provides a useful additional degree of freedom. Our work in this paper stands at the intersection of these two approaches. To our knowledge, that intersection was so far empty: no meaningful closed-form formulation was known for regularized optimal transport. We provide closed-form formulas of entropic (OT) of two Gaussian measures for balanced and unbalanced cases.

Summarizing measures vs. regularizing OT. Closed-form identities to compute OT distances (or more generally recover Monge maps) are known when either (1) both measures are univariate and the ground cost is submodular [44, §2]: in that case evaluating OT only requires integrating that submodular cost w.r.t. the quantile distributions of both measures; or (2) both measures are Gaussian, in a Hilbert space, and the ground cost is the squared Euclidean metric [19, 24], in which case the OT cost is given by the Wasserstein-Bures metric [5, 36]. These two formulas have inspired several works in which data measures are either projected onto 1D lines [42, 7], with further developments in [40, 32, 48]; or represented by Gaussians, to take advantage of the simpler computational possibilities offered by the Wasserstein-Bures metric [29, 39, 12].

Various schemes have been proposed to regularize the OT problem in the primal [15, 23] or the dual [46, 2, 16]. We focus in this work on the formulation obtained by [14], which combines entropic regularization [15] with a more general formulation for unbalanced transport [13, 33, 34]. The advantages of unbalanced entropic transport are numerous: it comes with favorable sample complexity regimes compared to unregularized OT [25], can be cast as a loss with favorable properties [27, 20], and can be evaluated using variations of the Sinkhorn algorithm [26].

On the absence of closed-form formulas for regularized OT. Despite its appeal, one of the shortcomings of entropic regularized OT lies in the absence of simple test-cases that admit closed-form formulas. While it is known that regularized OT can be related, in the limit of infinite regularization, to the energy distance [43], the absence of closed-form formulas for a fixed regularization strength poses an important practical problem to evaluate the performance of stochastic algorithms that try to approximate regularized OT: we do not know of any setup for which the ground truth value of entropic OT between continuous densities is known. The purpose of this paper is to fill this gap, and provide closed form expressions for balanced and unbalanced OT for Gaussian measures. We hope these formulas will prove useful in two different ways: as a solution to the problem outlined above, to facilitate the evaluation of new methodologies building on entropic OT, and more generally to propose a more robust yet well-grounded replacement to the Bures-Wasserstein metric.

Related work. From an economics theory perspective, Bojilov and Galichon [6] provided a closed form for an “equilibrium 2-sided matching problem” which is equivalent to entropy-regularized optimal transport. Second, a sequence of works in optimal control theory [10, 11, 9] studied stochastic systems, of which entropy regularized optimal transport between Gaussians can be seen as a special case, and found a closed form of the optimal dual potentials. Finally, a few recent concurrent works provided a closed form of entropy regularized OT between Gaussians: first Gerolin et al. [28] found a closed form in the univariate case, then Mallasto et al. [37] and del Barrio and Loubes [17] generalized the formula for multivariate Gaussians. The closest works to this paper are certainly those of Mallasto et al. [37] and del Barrio and Loubes [17] where the authors solved the balanced entropy regularized OT and studied the Gaussian barycenters problem. To the best of our knowledge, the closed form formula we provide for unbalanced OT is novel. Other differences between this paper and the aforementioned papers are highlighted below.

Contributions. Our contributions can be summarized as follows:

- Theorem 1 provides a closed form expression of the entropic (OT) plan π , which is shown to be a Gaussian measure itself (also shown in [6, 9, 37, 17]). Here, we furthermore study the properties of the OT loss function: it remains well defined, convex and differentiable even for singular covariance matrices unlike the Bures metric.

- Using the definition of debiased Sinkhorn barycenters [35, 31], Theorem 2 shows that the entropic barycenter of Gaussians is Gaussian and its covariance verifies a fixed point equation similar to that of Agueh and Carlier [1]. Mallasto et al. [37] and del Barrio and Loubes [17] provided similar fix point equations however by restricting the barycenter problem to the set of Gaussian measures whereas we consider the larger set of sub-Gaussian measures.
- As in the balanced case, Theorem 3 provides a closed form expression of the unbalanced Gaussian transport plan. The obtained formula sheds some light on the link between mass destruction and the distance between the means of α, β in Unbalanced OT.

Notations. \mathcal{S}^d denotes the set of square symmetric matrices in $\mathbb{R}^{d \times d}$. \mathcal{S}_{++}^d and \mathcal{S}_+^d denote the cones of positive definite and positive semi-definite matrices in \mathcal{S}^d respectively. Let $\mathcal{N}(\mathbf{a}, \mathbf{A})$ denote the multivariate Gaussian distribution with mean $\mathbf{a} \in \mathbb{R}^d$ and variance $\mathbf{A} \in \mathcal{S}_{++}^d$. $f = \mathcal{Q}(\mathbf{a}, \mathbf{A})$ denotes the quadratic form $f : x \mapsto -\frac{1}{2}(x^\top \mathbf{A} x - 2\mathbf{a}^\top x)$ with $\mathbf{A} \in \mathcal{S}^d$. For short, we denote $\mathcal{Q}(\mathbf{A}) = \mathcal{Q}(0, \mathbf{A})$. Whenever relevant, we follow the convention $0 \log 0 = 0$. \mathcal{M}_p^+ denotes the set of non-negative measures in \mathbb{R}^d with a finite p-th order moment and its subset of probability measures \mathcal{P}_p . For a non-negative measure $\alpha \in \mathcal{M}_p^+(\mathbb{R}^d)$, $\mathcal{L}_2(\alpha)$ denotes the set of functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\mathbb{E}_\alpha(|f|^2) = \int_{\mathbb{R}^d} |f|^2 d\alpha < +\infty$. With $\mathbf{C} \in \mathcal{S}_{++}^d$ and $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$, we denote the squared Mahalanobis distance: $\|\mathbf{a} - \mathbf{b}\|_{\mathbf{C}}^2 = (\mathbf{a} - \mathbf{b})^\top \mathbf{C} (\mathbf{a} - \mathbf{b})$.

2 Reminders on Optimal Transport

The Kantorovich problem. Let $\alpha, \beta \in \mathcal{P}_2$ and let $\Pi(\alpha, \beta)$ denote the set of probability measures in \mathcal{P}_2 with marginal distributions equal to α and β . The 2-Wasserstein distance is defined as:

$$W_2^2(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y). \quad (1)$$

This is known as the *Kantorovich* formulation of optimal transport. When α is absolutely continuous with respect to the Lebesgue measure (i.e. when α has a density), Equation (1) can be equivalently rewritten using the *Monge* formulation, where $T_\# \mu = \nu$ i.f.f. for all Borel sets A , $\nu(T(A)) = \mu(A)$:

$$W_2^2(\alpha, \beta) = \min_{T: T_\# \alpha = \nu} \int_{\mathbb{R}^d} \|x - T(x)\|^2 d\alpha(x). \quad (2)$$

The optimal map T^* in Equation (2) is called the Monge map.

The Wasserstein-Bures metric. Let $\mathcal{N}(m, \Sigma)$ denote the Gaussian distribution on \mathbb{R}^d with mean $m \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathcal{S}_{++}^d$. A well-known fact [19, 47] is that Equation (1) admits a closed form for Gaussian distributions, called the Wasserstein-Bures distance (a.k.a. the *Fréchet* distance):

$$W_2^2(\mathcal{N}(a, \mathbf{A}), \mathcal{N}(b, \mathbf{B})) = \|a - b\|^2 + \mathfrak{B}^2(\mathbf{A}, \mathbf{B}), \quad (3)$$

where \mathfrak{B} is the *Bures* distance [5] between positive matrices:

$$\mathfrak{B}^2(\mathbf{A}, \mathbf{B}) \stackrel{\text{def}}{=} \text{Tr} \mathbf{A} + \text{Tr} \mathbf{B} - 2\text{Tr}(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}. \quad (4)$$

Moreover, the Monge map between two Gaussian distributions admits a closed form: $T^* : x \rightarrow \mathbf{T}^{\mathbf{AB}}(x - \mathbf{a}) + \mathbf{b}$, with

$$\mathbf{T}^{\mathbf{AB}} \stackrel{\text{def}}{=} \mathbf{A}^{-\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}} = \mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}})^{-\frac{1}{2}} \mathbf{B}^{\frac{1}{2}}, \quad (5)$$

which is related to the Bures gradient (w.r.t. the Frobenius inner product):

$$\nabla_{\mathbf{A}} \mathfrak{B}^2(\mathbf{A}, \mathbf{B}) = \text{Id} - \mathbf{T}^{\mathbf{AB}}. \quad (6)$$

$\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$ and its gradient can be computed efficiently on GPUs using Newton-Schulz iterations which are provided in Algorithm 1 along with numerical experiments in the appendix.

3 Entropy-Regularized Optimal Transport between Gaussians

Solving (1) can be quite challenging, even in a discrete setting [41]. Adding an entropic regularization term to (1) results in a problem which can be solved efficiently using Sinkhorn's algorithm [15]. Let $\sigma > 0$. This corresponds to solving the following problem:

$$\text{OT}_\sigma(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y) + 2\sigma^2 \text{KL}(\pi \| \alpha \otimes \beta), \quad (7)$$

where $\text{KL}(\pi \| \alpha \otimes \beta) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} \log \left(\frac{d\pi}{d\alpha d\beta} \right) d\pi$ is the Kullback-Leibler divergence (or relative entropy). As in the original case (1), OT_σ can be studied with centered measures (i.e zero mean) with no loss of generality:

Lemma 1. *Let $\alpha, \beta \in \mathcal{P}$ and $\bar{\alpha}, \bar{\beta}$ their respective centered transformations. It holds that*

$$\text{OT}_\sigma(\alpha, \beta) = \text{OT}_\sigma(\bar{\alpha}, \bar{\beta}) + \|\mathbf{a} - \mathbf{b}\|^2. \quad (8)$$

Dual problem and Sinkhorn's algorithm. Compared to (1), (7) enjoys additional properties, such as the uniqueness of the solution π^* . Moreover, problem (7) has the following dual formulation:

$$\text{OT}_\sigma(\alpha, \beta) = \max_{\substack{f \in \mathcal{L}_1(\alpha), \\ g \in \mathcal{L}_1(\beta)}} \mathbb{E}_\alpha(f) + \mathbb{E}_\beta(g) - 2\sigma^2 \left(\int_{\mathbb{R}^d \times \mathbb{R}^d} e^{\frac{f(x) + g(y) - \|x - y\|^2}{2\sigma^2}} d\alpha(x) d\beta(y) - 1 \right). \quad (9)$$

If α and β have finite second order moments, a pair of dual potentials (f, g) is optimal if and only they verify the following optimality conditions β -a.s and α -a.s respectively [38]:

$$e^{\frac{f(x)}{2\sigma^2}} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x - y\|^2 + g(y)}{2\sigma^2}} d\beta(y) \right) = 1, \quad e^{\frac{g(y)}{2\sigma^2}} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x - y\|^2 + f(x)}{2\sigma^2}} d\alpha(x) \right) = 1. \quad (10)$$

Moreover, given a pair of optimal dual potentials (f, g) , the optimal transportation plan is given by

$$\frac{d\pi^*}{d\alpha d\beta}(x, y) = e^{\frac{f(x) + g(y) - \|x - y\|^2}{2\sigma^2}}. \quad (11)$$

Starting from a pair of potentials (f_0, g_0) , the optimality conditions (10) lead to an alternating dual ascent algorithm, which is equivalent to Sinkhorn's algorithm in log-domain:

$$\begin{aligned} g_{n+1} &= \left(y \in \mathbb{R}^d \rightarrow -2\sigma^2 \log \int_{\mathbb{R}^d} e^{\frac{-\|x - y\|^2 + f_n(x)}{2\sigma^2}} d\alpha(x) \right), \\ f_{n+1} &= \left(x \in \mathbb{R}^d \rightarrow -2\sigma^2 \log \int_{\mathbb{R}^d} e^{\frac{-\|x - y\|^2 + g_{n+1}(y)}{2\sigma^2}} d\beta(y) \right). \end{aligned} \quad (12)$$

Séjourné et al. [45] showed that when the support of the measures is compact, Sinkhorn's algorithm converges to a pair of dual potentials. Here in particular, we study Sinkhorn's algorithm when α and β are Gaussian measures.

Closed form expression for Gaussian measures.

Theorem 1. *Let $\mathbf{A}, \mathbf{B} \in \mathcal{S}_{++}^d$ and $\alpha \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta \sim \mathcal{N}(\mathbf{b}, \mathbf{B})$. Define $\mathbf{D}_\sigma = (4\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}} + \sigma^4 \text{Id})^{\frac{1}{2}}$. Then,*

$$\text{OT}_\sigma(\alpha, \beta) = \|\mathbf{a} - \mathbf{b}\|^2 + \mathcal{B}_\sigma^2(\mathbf{A}, \mathbf{B}), \quad \text{where} \quad (13)$$

$$\mathcal{B}_\sigma^2(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - \text{Tr}(\mathbf{D}_\sigma) + d\sigma^2(1 - \log(2\sigma^2)) + \sigma^2 \log \det(\mathbf{D}_\sigma + \sigma^2 \text{Id}). \quad (14)$$

Moreover, with $\mathbf{C}_\sigma = \frac{1}{2}\mathbf{A}^{\frac{1}{2}}\mathbf{D}_\sigma\mathbf{A}^{-\frac{1}{2}} - \frac{\sigma^2}{2}\text{Id}$, the Sinkhorn optimal transportation plan is also a Gaussian measure over $\mathbb{R}^d \times \mathbb{R}^d$ given by

$$\pi^* \sim \mathcal{N}\left(\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}, \begin{pmatrix} \mathbf{A} & \mathbf{C}_\sigma \\ \mathbf{C}_\sigma^\top & \mathbf{B} \end{pmatrix}\right). \quad (15)$$

Remark 1. While for our proof it is necessary to assume that \mathbf{A} and \mathbf{B} are positive definite in order for them to have a Lebesgue density, notice that the closed form formula given by Theorem 1 remains well-defined for positive semi-definite matrices. Moreover, unlike the Bures-Wasserstein metric, OT_σ is differentiable even when \mathbf{A} or \mathbf{B} are singular.

The proof of 1 is broken down into smaller results, Propositions 1 to 3 and lemma 2. Using Lemma 1, we can focus in the rest of this section on centered Gaussians without loss of generality.

Sinkhorn's algorithm and quadratic potentials. We obtain a closed form solution of OT_σ by considering quadratic solutions of (10). The following key proposition characterizes the obtained potential after a pair of Sinkhorn iterations with quadratic forms.

Proposition 1. *Let $\alpha \sim \mathcal{N}(0, \mathbf{A})$ and $\beta \sim \mathcal{N}(0, \mathbf{B})$ and the Sinkhorn transform $T_\alpha : \mathbb{R}^{\mathbb{R}^d} \rightarrow \mathbb{R}^{\mathbb{R}^d}$:*

$$T_\alpha(h)(x) \stackrel{\text{def}}{=} -\log \int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2}{2\sigma^2} + h(y)} d\alpha(y). \quad (16)$$

Let $\mathbf{X} \in \mathcal{S}_d$. If $h = m + \mathcal{Q}(\mathbf{X})$ i.e $h(x) = m - \frac{1}{2}x^\top \mathbf{X}x$ for some $m \in \mathbb{R}$, then $T_\alpha(h)$ is well-defined if and only if $\mathbf{X}' \stackrel{\text{def}}{=} \sigma^2 \mathbf{X} + \sigma^2 \mathbf{A}^{-1} + \text{Id} \succ 0$. In that case,

- (i) $T_\alpha(h) = \mathcal{Q}(\mathbf{Y}) + m'$ where $\mathbf{Y} = \frac{1}{\sigma^2}(\mathbf{X}'^{-1} - \text{Id})$ and $m' \in \mathbb{R}$ is an additive constant,
- (ii) $T_\beta(T_\alpha(h))$ is well-defined and is also a quadratic form up to an additive constant, since $\mathbf{Y}' \stackrel{\text{def}}{=} \sigma^2 \mathbf{Y} + \sigma^2 \mathbf{B}^{-1} + \text{Id} = \mathbf{X}'^{-1} + \sigma^2 \mathbf{B}^{-1} \succ 0$ and (i) applies.

Consider the null inialization $f_0 = 0 = \mathcal{Q}(0)$. Since $\sigma^2 \mathbf{A}^{-1} + \text{Id} \succ 0$, Proposition 1 applies with $\mathbf{X} = 0$ and a simple induction shows that (f_n, g_n) remain quadratic forms for all n . Sinkhorn's algorithm can thus be written as an algorithm on positive definite matrices.

Proposition 2. *Starting with null potentials, Sinkhorn's algorithm is equivalent to the iterations:*

$$\mathbf{F}_{n+1} = \sigma^2 \mathbf{A}^{-1} + \mathbf{G}_n^{-1}, \quad \mathbf{G}_{n+1} = \sigma^2 \mathbf{B}^{-1} + \mathbf{F}_{n+1}^{-1}, \quad (17)$$

with $\mathbf{F}_0 = \sigma^2 \mathbf{A}^{-1} + \text{Id}$ and $\mathbf{G}_0 = \sigma^2 \mathbf{B}^{-1} + \text{Id}$.

Moreover, the sequence $(\mathbf{F}_n, \mathbf{G}_n)$ is contractive (in the matrix operator norm) and converges towards a pair of positive definite matrices (\mathbf{F}, \mathbf{G}) . At optimality, the dual potentials are determined up to additive constants f_0 and g_0 : $\frac{f}{2\sigma^2} = \mathcal{Q}(\mathbf{U}) + f_0$ and $\frac{g}{2\sigma^2} = \mathcal{Q}(\mathbf{V}) + g_0$ where \mathbf{U} and \mathbf{V} are given by

$$\mathbf{F} = \sigma^2 \mathbf{U} + \sigma^2 \mathbf{A}^{-1} + \text{Id}, \quad \mathbf{G} = \sigma^2 \mathbf{V} + \sigma^2 \mathbf{B}^{-1} + \text{Id}. \quad (18)$$

Closed form solution. Taking the limit of Sinkhorn's equations (17) along with the change of variable (18), there exists a pair of optimal potentials determined up to an additive constant:

$$\frac{f}{2\sigma^2} = \mathcal{Q}(\mathbf{U}) = \mathcal{Q}\left(\frac{1}{\sigma^2}(\mathbf{G}^{-1} - \text{Id})\right), \quad \frac{g}{2\sigma^2} = \mathcal{Q}(\mathbf{V}) = \mathcal{Q}\left(\frac{1}{\sigma^2}(\mathbf{F}^{-1} - \text{Id})\right), \quad (19)$$

where (\mathbf{F}, \mathbf{G}) is the solution of the fixed point equations

$$\mathbf{F} = \sigma^2 \mathbf{A}^{-1} + \mathbf{G}^{-1}, \quad \mathbf{G} = \sigma^2 \mathbf{B}^{-1} + \mathbf{F}^{-1}. \quad (20)$$

Let $\mathbf{C} \stackrel{\text{def}}{=} \mathbf{A}\mathbf{G}^{-1}$. Combining both equations of (20) in one leads to $\mathbf{G} = \sigma^2 \mathbf{B}^{-1} + (\mathbf{G}^{-1} + \sigma^2 \mathbf{A}^{-1})^{-1}$, which can be shown to be equivalent to

$$\mathbf{C}^2 + \sigma^2 \mathbf{C} - \mathbf{A}\mathbf{B} = 0. \quad (21)$$

Notice that since \mathbf{A} and \mathbf{G}^{-1} are positive definite, their product $\mathbf{C} = \mathbf{A}\mathbf{G}^{-1}$ is similar to $\mathbf{A}^{\frac{1}{2}}\mathbf{G}^{-1}\mathbf{A}^{\frac{1}{2}}$. Thus it has positive eigenvalues. Proposition 3 provides the only feasible solution of (21).

Proposition 3. *Let $\sigma^2 \geq 0$ and \mathbf{C} satisfying Equation (21). Then,*

$$\mathbf{C} = \left(\mathbf{A}\mathbf{B} + \frac{\sigma^4}{4} \text{Id}\right)^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id} = \mathbf{A}^{\frac{1}{2}}(\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}}\mathbf{A}^{-\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id}. \quad (22)$$

Corollary 1. *The optimal dual potentials of (19) can be given in closed form by:*

$$\mathbf{U} = \frac{\mathbf{B}}{\sigma^2}(\mathbf{C} + \sigma^2 \text{Id})^{-1} - \frac{\text{Id}}{\sigma^2}, \quad \mathbf{V} = (\mathbf{C} + \sigma^2 \text{Id})^{-1} \frac{\mathbf{A}}{\sigma^2} - \frac{\text{Id}}{\sigma^2}. \quad (23)$$

Moreover, \mathbf{U} and \mathbf{V} remain well-defined even for singular matrices \mathbf{A} and \mathbf{B} .

Optimal transportation plan and OT_σ . Using Corollary 1 and (19), Equation (11) leads to a closed form expression of π . To conclude the proof of Theorem 1, we introduce lemma 2 that computes the OT_σ loss at optimality. Detailed technical proofs are provided in the appendix.

Lemma 2. *Let $\mathbf{A}, \mathbf{B}, \mathbf{C}$ be invertible matrices such that $\mathbf{H} = \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \succ 0$. Let $\alpha = \mathcal{N}(0, \mathbf{A}), \beta = \mathcal{N}(0, \mathbf{B})$, and $\pi = \mathcal{N}(0, \mathbf{H})$. Then,*

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - 2\text{Tr}(\mathbf{C}), \quad (24)$$

$$\text{KL}(\pi \| \alpha \otimes \beta) = \frac{1}{2} (\log \det \mathbf{A} + \log \det \mathbf{B} - \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}). \quad (25)$$

Properties of OT_σ . Theorem 1 shows that π has a Gaussian density. Proposition 4 allows to reformulate this optimization problem over couplings in $\mathbb{R}^{d \times d}$ with a positivity constraint.

Proposition 4. *Let $\alpha = \mathcal{N}(0, \mathbf{A}), \beta = \mathcal{N}(0, \mathbf{B})$, and $\sigma^2 > 0$. Then,*

$$\text{OT}_\sigma(\alpha, \beta) = \min_{\mathbf{C}: \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \geq 0} \left\{ \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - 2\text{Tr}(\mathbf{C}) + \sigma^2 (\log \det \mathbf{A} \mathbf{B} - \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}) \right\} \quad (26)$$

$$= \min_{\mathbf{K} \in \mathbb{R}^{d \times d}; \|\mathbf{K}\|_{op} \leq 1} \text{Tr} \mathbf{A} + \text{Tr} \mathbf{B} - 2\text{Tr} \mathbf{A}^{\frac{1}{2}} \mathbf{K} \mathbf{B}^{\frac{1}{2}} - \sigma^2 \ln \det(\text{Id} - \mathbf{K} \mathbf{K}^\top). \quad (27)$$

Moreover, both (26) and (27) are convex problems.

We now study the convexity and differentiability of OT_σ , which are more conveniently derived from the dual problem of (26) given as a positive definite program:

Proposition 5. *The dual problem of (26) can be written with no duality gap as*

$$\max_{\mathbf{F}, \mathbf{G} \succ 0} \left\{ \langle \text{Id} - \mathbf{F}, \mathbf{A} \rangle + \langle \text{Id} - \mathbf{G}, \mathbf{B} \rangle + \sigma^2 \log \det \left(\frac{\mathbf{F} \mathbf{G} - \text{Id}}{\sigma^4} \right) + \sigma^2 \log \det \mathbf{A} \mathbf{B} + 2d\sigma^2 \right\}. \quad (28)$$

Feydy et al. [20] showed that on compact spaces, the gradient of OT_σ is given by the optimal dual potentials. This result was later extended by Janati et al. [31] to sub-Gaussian measures with unbounded supports. The following proposition re-establishes this statement for Gaussians.

Proposition 6. *Assume $\sigma > 0$ and consider the pair \mathbf{U}, \mathbf{V} of Corollary 1. Then*

- (i) *The optimal pair $(\mathbf{F}^*, \mathbf{G}^*)$ of (28) is a solution to the fixed point problem (20),*
- (ii) *\mathfrak{B}_{σ^2} is differentiable and: $\nabla \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B}) = -(\sigma^2 \mathbf{U}, \sigma^2 \mathbf{V})$. Thus: $\nabla_{\mathbf{A}} \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B}) = \text{Id} - \mathbf{B}^{\frac{1}{2}} \left((\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}}$,*
- (iii) *$(\mathbf{A}, \mathbf{B}) \mapsto \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$ is convex in \mathbf{A} and in \mathbf{B} but not jointly.*
- (iv) *For a fixed \mathbf{B} with its spectral decomposition $\mathbf{B} = \mathbf{P} \Sigma \mathbf{P}^\top$, the function $\phi_{\mathbf{B}} : \mathbf{A} \mapsto \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$ is minimized at $\mathbf{A}_0 = \mathbf{P}(\Sigma - \sigma^2 \text{Id})_+ \mathbf{P}^\top$ where the thresholding operator $_+$ is defined by $x_+ = \max(x, 0)$ for any $x \in \mathbb{R}$ and extended element-wise to diagonal matrices.*

When \mathbf{A} and \mathbf{B} are not singular, by letting $\sigma \rightarrow 0$ in $\nabla_{\mathbf{A}} \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$, we recover the gradient of the Bures metric given in (6). Moreover, (iv) illustrates the entropy bias of \mathfrak{B}_{σ^2} . Feydy et al. [20] showed that it can be circumvented by considering the Sinkhorn divergence:

$$S_\sigma : (\alpha, \beta) \mapsto \text{OT}_\sigma(\alpha, \beta) - \frac{1}{2} (\text{OT}_\sigma(\alpha, \alpha) + \text{OT}_\sigma(\beta, \beta)) \quad (29)$$

which is non-negative and equals 0 if and only if $\alpha = \beta$. Using the differentiability and convexity of S_σ on sub-Gaussian measures [31], we conclude this section by showing that the debiased Sinkhorn barycenter of Gaussians remains Gaussian:

Theorem 2. *Consider the restriction of OT_σ to the set of sub-Gaussian measures $\mathcal{G} \stackrel{\text{def}}{=} \{\mu \in \mathcal{P}_2 | \exists q > 0, \mathbb{E}_\mu(e^{q\|X\|^2}) < +\infty\}$ and let K Gaussian measures $\alpha_k \sim \mathcal{N}(\mathbf{a}_k, \mathbf{A}_k)$ with a sequence of positive weights $(w_k)_k$ summing to 1. Then, the weighted debiased barycenter defined by:*

$$\beta \stackrel{\text{def}}{=} \argmin_{\beta \in \mathcal{G}} \sum_{k=1} w_k S_\sigma(\alpha_k, \beta) \quad (30)$$

is a Gaussian measure given by $\mathcal{N}\left(\sum_{k=1}^K w_k \mathbf{a}_k, \mathbf{B}\right)$ where $\mathbf{B} \in \mathcal{S}_+^d$ is a solution of the equation:

$$\sum_{k=1}^K w_k (\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} = (\mathbf{B}^2 + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} \quad (31)$$

4 Entropy Regularized OT between Unbalanced Gaussians

We proceed by considering a more general setting, in which measures $\alpha, \beta \in \mathcal{M}_2^+(\mathbb{R}^d)$ have finite integration masses $m_\alpha = \alpha(\mathbb{R}^d)$ and $m_\beta = \beta(\mathbb{R}^d)$ that are not necessarily the same. Following [14], we define entropy-regularized unbalanced OT as:

$$\text{UOT}_\sigma(\alpha, \beta) \stackrel{\text{def}}{=} \inf_{\pi \in \mathcal{M}_2^+(\mathbb{R}^d \times \mathbb{R}^d)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x-y\|^2 d\pi(x, y) + 2\sigma^2 \text{KL}(\pi \| \alpha \otimes \beta) + \gamma \text{KL}(\pi_1 \| \alpha) + \gamma \text{KL}(\pi_2 \| \beta), \quad (32)$$

where $\gamma > 0$ and π_1, π_2 are the marginal distributions of the coupling $\pi \in \mathcal{M}_2^+(\mathbb{R}^2 \times \mathbb{R}^d)$.

Duality and optimality conditions. By definition of the KL divergence, the term $\text{KL}(\pi \| \alpha \otimes \beta)$ in (32) is finite if and only if π admits a density with respect to $\alpha \otimes \beta$. Therefore (32) can be formulated as a variational problem:

$$\begin{aligned} \text{UOT}_\sigma(\alpha, \beta) \stackrel{\text{def}}{=} \inf_{r \in \mathcal{L}_1(\alpha \otimes \beta)} \left\{ \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x-y\|^2 r(x, y) d\alpha(x) d\beta(y) \right. \\ \left. + 2\sigma^2 \text{KL}(r \| \alpha \otimes \beta) + \gamma \text{KL}(r_1 \| \alpha) + \gamma \text{KL}(r_2 \| \beta) \right\}, \end{aligned} \quad (33)$$

where $r_1 \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} r(\cdot, y) d\beta(y)$ and $r_2 \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} r(x, \cdot) d\alpha(x)$ correspond to the marginal density functions and the Kullback-Leibler divergence is defined as: $\text{KL}(f \| \mu) = \int_{\mathbb{R}^d} (f \log(f) + f - 1) d\mu$. As in [14], Fenchel-Rockafellar duality holds and (33) admits the following dual problem:

$$\begin{aligned} \text{UOT}_\sigma(\alpha, \beta) = \sup_{\substack{f \in \mathcal{L}_\infty(\alpha) \\ g \in \mathcal{L}_\infty(\beta)}} \left\{ \gamma \int_{\mathbb{R}^d} (1 - e^{-\frac{f}{\gamma}}) d\alpha + \gamma \int_{\mathbb{R}^d} (1 - e^{-\frac{g}{\gamma}}) d\beta \right. \\ \left. - 2\sigma^2 \int_{\mathbb{R}^d \times \mathbb{R}^d} (e^{\frac{-\|x-y\|^2 + f(x) + g(y)}{2\sigma^2}} - 1) d\alpha(x) d\beta(y) \right\}, \end{aligned} \quad (34)$$

for which the necessary optimality conditions read, with $\tau \stackrel{\text{def}}{=} \frac{\gamma}{\gamma + 2\sigma^2}$:

$$\frac{f(x)}{2\sigma^2} \stackrel{a.s.}{=} -\tau \log \int_{\mathbb{R}^d} e^{\frac{g(y) - \|x-y\|^2}{2\sigma^2}} d\beta(y), \quad \frac{g(y)}{2\sigma^2} \stackrel{a.s.}{=} -\tau \log \int_{\mathbb{R}^d} e^{\frac{f(x) - \|x-y\|^2}{2\sigma^2}} d\alpha(x). \quad (35)$$

Moreover, if such a pair of dual potentials exists, then the optimal transportation plan is given by

$$\frac{d\pi}{d\alpha \otimes d\beta}(x, y) = e^{\frac{f(x) + g(y) - \|x-y\|^2}{2\sigma^2}}. \quad (36)$$

The following proposition provides a simple formula to compute UOT_σ at optimality. It shows that it is sufficient to know the total transported mass $\pi(\mathbb{R}^d \times \mathbb{R}^d)$.

Proposition 7. Assume there exists an optimal transportation plan π^* , solution of (32). Then

$$\text{UOT}_\sigma(\alpha, \beta) = \gamma(m_\alpha + m_\beta) + 2\sigma^2 m_\alpha m_\beta - 2(\sigma^2 + \gamma) \pi^*(\mathbb{R}^d \times \mathbb{R}^d). \quad (37)$$

Unbalanced OT for scaled Gaussians. Let α and β be unbalanced Gaussian measures. Formally, $\alpha = m_\alpha \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta = m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ with $m_\alpha, m_\beta > 0$. Unlike balanced OT, α and β cannot be assumed to be centered without loss of generality. However, we can still derive a closed form formula for $\text{UOT}_\sigma(\alpha, \beta)$ by considering quadratic potentials of the form

$$\frac{f(\mathbf{x})}{2\sigma^2} = -\frac{1}{2}(\mathbf{x}^\top \mathbf{U} \mathbf{x} - 2\mathbf{x}^\top \mathbf{u}) + \log(m_u), \quad \frac{g(\mathbf{y})}{2\sigma^2} = -\frac{1}{2}(\mathbf{x}^\top \mathbf{V} \mathbf{x} - 2\mathbf{x}^\top \mathbf{v}) + \log(m_v). \quad (38)$$

Let σ and γ be the regularization parameters as in Equation (33), and $\tau \stackrel{\text{def}}{=} \frac{\gamma}{2\sigma^2 + \gamma}$, $\lambda \stackrel{\text{def}}{=} \frac{\sigma^2}{1 - \tau} = \sigma^2 + \frac{\gamma}{2}$. Let us define the following useful quantities:

$$\mu = \begin{pmatrix} \mathbf{a} + \mathbf{A}\mathbf{X}^{-1}(\mathbf{b} - \mathbf{a}) \\ \mathbf{b} + \mathbf{B}\mathbf{X}^{-1}(\mathbf{a} - \mathbf{b}) \end{pmatrix} \quad (39)$$

$$\mathbf{H} = \begin{pmatrix} (\text{Id} + \frac{1}{\lambda}\mathbf{C})(\mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}) & \mathbf{C} + (\text{Id} + \frac{1}{\lambda}\mathbf{C})\mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{C}^\top + (\text{Id} + \frac{1}{\lambda}\mathbf{C}^\top)\mathbf{B}\mathbf{X}^{-1}\mathbf{A} & (\text{Id} + \frac{1}{\lambda}\mathbf{C}^\top)(\mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B}) \end{pmatrix} \quad (40)$$

$$m_\pi = \sigma^{\frac{d\sigma^2}{\gamma + \sigma^2}} \left(m_\alpha m_\beta \det(\mathbf{C}) \sqrt{\frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})^\tau}{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})}} \right)^{\frac{1}{\tau+1}} \frac{e^{-\frac{\|\mathbf{a} - \mathbf{b}\|_{\tilde{\mathbf{X}}}^2}{2(\tau+1)}}}{\sqrt{\det(\mathbf{C} - \frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}})}}, \quad (41)$$

with

$$\begin{aligned} \mathbf{X} &= \mathbf{A} + \mathbf{B} + \lambda \text{Id}, & \tilde{\mathbf{A}} &= \frac{\gamma}{2}(\text{Id} - \lambda(\mathbf{A} + \lambda \text{Id})^{-1}), \\ \tilde{\mathbf{B}} &= \frac{\gamma}{2}(\text{Id} - \lambda(\mathbf{B} + \lambda \text{Id})^{-1}), & \mathbf{C} &= \left(\frac{1}{\tau}\tilde{\mathbf{A}}\tilde{\mathbf{B}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id}. \end{aligned}$$

Theorem 3. Let $\alpha = m_\alpha \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta = m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ be two unbalanced Gaussian measures. Let $\tau = \frac{\gamma}{2\sigma^2 + \gamma}$ and $\lambda \stackrel{\text{def}}{=} \frac{\sigma^2}{1 - \tau} = \sigma^2 + \frac{\gamma}{2}$ and μ, \mathbf{H} , and m_π be as above. Then

- (i) The unbalanced optimal transport plan, minimizer of (32), is also an unbalanced Gaussian over $\mathbb{R}^d \times \mathbb{R}^d$ given by $\pi = m_\pi \mathcal{N}(\mu, \mathbf{H})$,
- (ii) UOT_σ can be obtained in closed form using Proposition 7 with $\pi(\mathbb{R}^d \times \mathbb{R}^d) = m_\pi$.

Remark 2. The exponential term in the closed form formula above provides some intuition on how transportation occurs in unbalanced OT. When the difference between the means is too large, the transported mass m_π^* goes to 0 and thus no transport occurs. However for fixed means \mathbf{a}, \mathbf{b} , when $\gamma \rightarrow +\infty$, $\mathbf{X}^{-1} \rightarrow 0$ and the exponential term approaches 1.

5 Numerical Experiments

Empirical validation of the closed form formulas. Figure 1 illustrates the convergence towards the closed form formulas of both theorems. For each dimension d in $[5, 10]$, we select a pair of Gaussians $\alpha \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta \sim m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ with m_β equals 1 (balanced) or 2 (unbalanced) and randomly generated means \mathbf{a}, \mathbf{b} (uniform in $[-1, 1]^d$) and covariances $\mathbf{A}, \mathbf{B} \in S_{++}^d$ following the Wishart distribution $W_d(0.2 * \text{Id}, d)$. We generate i.i.d datasets $\alpha_n \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta_n \sim m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ with n samples and compute $\text{OT}_\sigma / \text{UOT}_\sigma$. We report means and \pm shaded standard-deviation areas over 20 independent trials for each value of n .

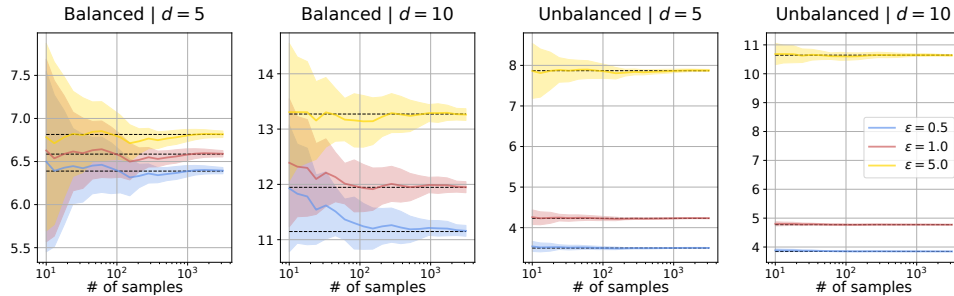


Figure 1: Numerical convergence the (n -samples) empirical estimation of $\text{OT}(\alpha_n, \beta_n)$ computed using Sinkhorn's algorithm towards the closed form of $\text{OT}_\sigma(\alpha, \beta)$ and $\text{UOT}_\sigma(\alpha, \beta)$ (the theoretical limit is dashed) given by Theorem 1 and Theorem 3 for random Gaussians α, β . For unbalanced OT, $\gamma = 1$.

Transport plan visualization with $d = 1$. Figure 2 confronts the expected theoretical plans (contours in black) given by theorems 1 and 3 to empirical ones (weights in shades of red) obtained

with Sinkhorn’s algorithm using 2000 Gaussian samples. The density functions (black) and the empirical histograms (red) of α (resp. β) with 200 bins are displayed on the left (resp. top) of each transport plan. The red weights are computed via a 2d histogram of the transport plan returned by Sinkhorn’s algorithm with (200×200) bins. Notice the blurring effect of ε and increased mass transportation of the Gaussian tails in unbalanced transport with larger γ .

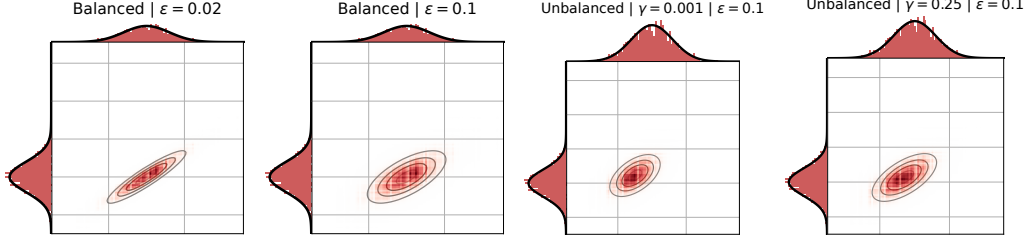


Figure 2: Effect of ε in balanced OT and γ in unbalanced OT. Empirical plans (red) correspond to the expected Gaussian contours depicted in black. Here $\alpha = \mathcal{N}(0, 0.04)$ and $\beta = m_\beta \mathcal{N}(0.5, 0.09)$ with $m_\beta = 1$ (balanced) and $m_\beta = 2$ (unbalanced). In unbalanced OT, the right tail of β is not transported, and the mean of the transportation plan is shifted compared to that of the balanced case – as expected from Theorem 3 specially for low γ .

Empirical estimation of the closed form mean and covariance of the unbalanced transport plan

Figure 3 illustrates the convergence towards the closed form formulas of μ and \mathbf{H} of theorem 3. For each dimension d in $[1, 2, 5, 10]$, we select a pair of Gaussians $\alpha \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta \sim m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ with $m_\beta = 1.1$ and randomly generated means \mathbf{a}, \mathbf{b} (uniform in $[-1, 1]^d$) and covariances $\mathbf{A}, \mathbf{B} \in S_{++}^d$ following the Wishart distribution $W_d(0.2 * \text{Id}, d)$. We generate i.i.d datasets $\alpha_n \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$ and $\beta_n \sim m_\beta \mathcal{N}(\mathbf{b}, \mathbf{B})$ with n samples and compute $\text{OT}_\sigma / \text{UOT}_\sigma$. We set $\varepsilon \stackrel{\text{def}}{=} 2\sigma^2 = 0.5$ and $\gamma = 0.1$. Using the obtained empirical Sinkhorn transportation plan, we computed its empirical mean μ_n and covariance matrix Σ_n and display their relative ℓ_∞ distance to μ and \mathbf{H} (Σ in the figure) of theorem 3. The means and \pm sd intervals are computed over 50 independent trials for each value of n .

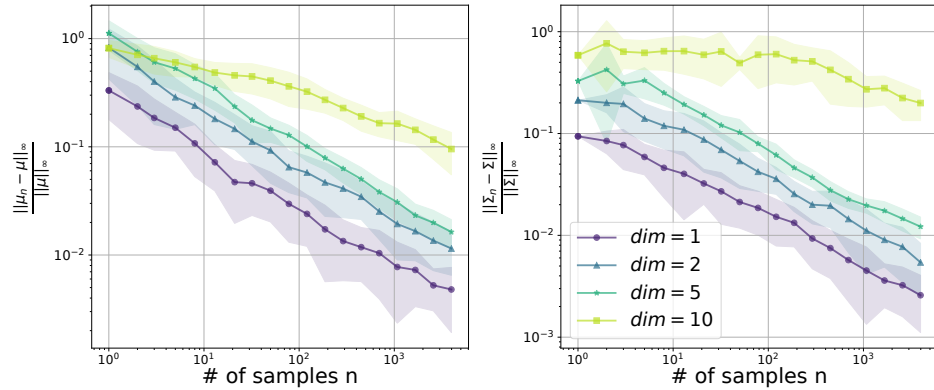


Figure 3: Numerical convergence the $(n\text{-samples})$ empirical estimation of the theoretical mean μ and covariance \mathbf{H} of theorem 3. Empirical moments are computed using Sinkhorn’s algorithm.

Broader Impact

We expect this work to benefit research on sample complexity issues in regularized optimal transport, such as [25] for balanced regularized OT, and future work on unbalanced regularized OT. By providing the first continuous test-case, we hope that researchers will be able to better test their theoretical bounds and benchmark their methods.

Acknowledgments

H. Janati, B. Muzellec and M. Cuturi were supported by a “Chaire d’excellence de l’IDEX Paris Saclay”. H. Janati acknowledges the support of the ERC Starting Grant SLAB ERC-YStG-676943. The work of G. Peyré was supported by the European Research Council (ERC project NORIA) and by the French government under management of Agence Nationale de la Recherche as part of the “Investissements d’avenir” program, reference ANR19-P3IA-0001 (PRAIRIE 3IA Institute).

References

- [1] Agueh, M. and Carlier, G. (2011). Barycenters in the Wasserstein space. *SIAM*, 43(2):904–924.
- [2] Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 214–223.
- [3] Benamou, J.-D. (2003). Numerical resolution of an “unbalanced” mass transport problem. *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, 37:851–868.
- [4] Bhatia, R. (2007). *Positive Definite Matrices*. Princeton Series in Applied Mathematics. Princeton University Press, Princeton, NJ, USA.
- [5] Bhatia, R., Jain, T., and Lim, Y. (2018). On the Bures-Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*.
- [6] Bojilov, R. and Galichon, A. (2016). Matching in closed-form: equilibrium, identification, and comparative statics. *Economic Theory*, 61(4):587–609.
- [7] Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. (2015). Sliced and Radon Wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45.
- [8] Bures, D. (1969). An extension of Kakutani’s theorem on infinite product measures to the tensor product of semifinite w^* -algebras. *Transactions of the American Mathematical Society*, 135:199–212.
- [9] Chen, Y., Georgiou, T. T., and Pavon, M. (2016). On the relation between optimal transport and Schrödinger bridges: A stochastic control viewpoint. *Journal of Optimization Theory and Applications*, 169(2):671–691.
- [10] Chen, Y., Georgiou, T. T., and Pavon, M. (2016). Optimal steering of a linear stochastic system to a final probability distribution, part i. *IEEE Transactions on Automatic Control*, 61(5):1158–1169.
- [11] Chen, Y., Georgiou, T. T., and Pavon, M. (2018). Optimal steering of a linear stochastic system to a final probability distribution—part iii. *IEEE Transactions on Automatic Control*, 63(9):3112–3118.
- [12] Chen, Y., Georgiou, T. T., and Tannenbaum, A. (2018). Optimal transport for Gaussian mixture models. *IEEE Access*, 7:6269–6278.
- [13] Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. (2018a). An interpolating distance between optimal transport and Fisher–Rao metrics. *Foundations of Computational Mathematics*, 18(1):1–44.
- [14] Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. (2018b). Scaling algorithms for unbalanced optimal transport problems. *Mathematics of Computation*, 87(314):2563–2609.
- [15] Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, pages 2292–2300.
- [16] Cuturi, M. and Peyré, G. (2016). A smoothed dual approach for variational Wasserstein problems. *SIAM Journal on Imaging Sciences*, 9(1):320–343.
- [17] del Barrio, E. and Loubes, J.-M. (2020). The statistical effect of entropic regularization in optimal transportation. *arxiv preprint arXiv:2006.05199*.
- [18] Dereich, S., Scheutzow, M., and Schottstedt, R. (2013). Constructive quantization: Approximation by empirical measures. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 49, pages 1183–1203.

- [19] Dowson, D. and Landau, B. (1982). The Fréchet distance between multivariate normal distributions. *Journal of Multivariate Analysis*, 12(3):450 – 455.
- [20] Feydy, J., Séjourné, T., Vialard, F., Amari, S., Trounev, A., and Peyré, G. (2019). Interpolating between optimal transport and MMD using Sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, pages 2681–2690.
- [21] Figalli, A. (2017). *The Monge–Ampère equation and its applications*.
- [22] Fournier, N. and Guillin, A. (2015). On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3-4):707–738.
- [23] Frogner, C., Zhang, C., Mobahi, H., Araya, M., and Poggio, T. A. (2015). Learning with a Wasserstein loss. In *Advances in Neural Information Processing Systems*, pages 2053–2061.
- [24] Gelbrich, M. (1990). On a formula for the l2 Wasserstein metric between measures on Euclidean and Hilbert spaces. *Mathematische Nachrichten*, 147(1):185–203.
- [25] Genevay, A., Chizat, L., Bach, F., Cuturi, M., and Peyré, G. (2019). Sample complexity of sinkhorn divergences. In Chaudhuri, K. and Sugiyama, M., editors, *Proceedings of Machine Learning Research*, volume 89 of *Proceedings of Machine Learning Research*, pages 1574–1583. PMLR.
- [26] Genevay, A., Cuturi, M., Peyré, G., and Bach, F. (2016). Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, pages 3440–3448.
- [27] Genevay, A., Peyre, G., and Cuturi, M. (2018). Learning generative models with Sinkhorn divergences. In Storkey, A. and Perez-Cruz, F., editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1608–1617, Playa Blanca, Lanzarote, Canary Islands. PMLR.
- [28] Gerolin, A., Grossi, J., and Gori-Giorgi, P. (2020). Kinetic correlation functionals from the entropic regularization of the strictly correlated electrons problem. *Journal of Chemical Theory and Computation*, 16(1):488–498.
- [29] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. (2017). Gans trained by a two time-scale update rule converge to a local Nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637.
- [30] Higham, N. J. (2008). *Functions of Matrices: Theory and Computation (Other Titles in Applied Mathematics)*. Society for Industrial and Applied Mathematics, USA.
- [31] Janati, H., Cuturi, M., and Gramfort, A. (2020). Debiased sinkhorn barycenters. In *Proceedings of the 34th International Conference on Machine Learning*.
- [32] Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., and Rohde, G. (2019). Generalized sliced wasserstein distances. In *Advances in Neural Information Processing Systems*, pages 261–272.
- [33] Liero, M., Mielke, A., and Savaré, G. (2016). Optimal transport in competition with reaction: the Hellinger–Kantorovich distance and geodesic curves. *SIAM Journal on Mathematical Analysis*, 48(4):2869–2911.
- [34] Liero, M., Mielke, A., and Savaré, G. (2018). Optimal entropy-transport problems and a new Hellinger–Kantorovich distance between positive measures. *Inventiones Mathematicae*, 211(3):969–1117.
- [35] Luise, G., Salzo, S., Pontil, M., and Ciliberto, C. (2019). Sinkhorn barycenters with free support via frank-wolfe algorithm. In *Advances in Neural Information Processing Systems*.
- [36] Malagò, L., Montrucchio, L., and Pistone, G. (2018). Wasserstein riemannian geometry of positive definite matrices. *arXiv preprint arXiv:1801.09269*.
- [37] Mallasto, A., Gerolin, A., and Minh, H. Q. (2020). Entropy-regularized 2-wasserstein distance between gaussian measures. *Arxiv preprint arXiv:2006.03416*.
- [38] Mena, G. and Niles-Weed, J. (2019). Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem. In *Advances in Neural Information Processing Systems* 32, pages 4541–4551. Curran Associates, Inc.
- [39] Muzellec, B. and Cuturi, M. (2018). Generalizing point embeddings using the wasserstein space of elliptical distributions. In *Advances in Neural Information Processing Systems* 31, pages 10237–10248. Curran Associates, Inc.

- [40] Paty, F.-P. and Cuturi, M. (2019). Subspace robust wasserstein distances. In *International Conference on Machine Learning*, pages 5072–5081.
- [41] Peyré, G. and Cuturi, M. (2019). Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–206.
- [42] Rabin, J., Peyré, G., Delon, J., and Bernot, M. (2011). Wasserstein barycenter and its application to texture mixing. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 435–446. Springer.
- [43] Ramdas, A., Trillos, N. G., and Cuturi, M. (2017). On wasserstein two-sample testing and related families of nonparametric tests. *Entropy*, 19(2):47.
- [44] Santambrogio, F. (2015). *Optimal transport for applied mathematicians*. Birkhauser.
- [45] Séjourné, T., Feydy, J., Vialard, F.-X., Trounev, A., and Peyré, G. (2019). Sinkhorn divergences for unbalanced optimal transport. *arXiv preprint arXiv:1910.12958*.
- [46] Shirdhonkar, S. and Jacobs, D. W. (2008). Approximate earth mover’s distance in linear time. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE.
- [47] Takatsu, A. (2011). Wasserstein geometry of Gaussian measures. *Osaka J. Math.*, 48(4):1005–1026.
- [48] Titouan, V., Flamary, R., Courty, N., Tavenard, R., and Chapel, L. (2019). Sliced gromov-wasserstein. In *Advances in Neural Information Processing Systems*, pages 14726–14736.
- [49] Villani, C. (2009). *Optimal Transport: Old and New*, volume 338. Springer Verlag.

Appendix

5.1 The Newton-Schulz algorithm

The main bottleneck in computing $\mathbf{T}^{\mathbf{AB}}$ is that of computing matrix square roots. This can be performed using singular value decomposition (SVD) or, as suggested in [39], using Newton-Schulz (NS) iterations [30, §5.3]. In particular, Newton-Schulz iterations have the advantage of yielding both roots, and inverse roots. Hence, to compute $\mathbf{T}^{\mathbf{AB}}$, one would run NS a first time to obtain $\mathbf{A}^{\frac{1}{2}}$ and $\mathbf{A}^{-\frac{1}{2}}$, and a second time to get $(\mathbf{A}^{\frac{1}{2}}\mathbf{B}\mathbf{A}^{\frac{1}{2}})^{\frac{1}{2}}$.

In fact, as a direct application of [30, Theorem 5.2], one can even compute both $\mathbf{T}^{\mathbf{AB}}$ and $\mathbf{T}^{\mathbf{BA}} = (\mathbf{T}^{\mathbf{AB}})^{-1}$ in a single run by initializing the Newton-Schulz algorithm with \mathbf{A} and \mathbf{B} , as in Algorithm 1. Using (6), and noting that $\mathfrak{B}^2(\mathbf{A}, \mathbf{B}) = \text{Tr}\mathbf{A} + \text{Tr}\mathbf{B} - 2\text{Tr}(\mathbf{T}^{\mathbf{AB}}\mathbf{A})$, this implies that a single run of NS is sufficient to compute $\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$, $\nabla_{\mathbf{A}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$ and $\nabla_{\mathbf{B}}\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$ using basic matrix operations. The main advantage of Newton-Schulz over SVD is that it its efficient scalability on GPUs, as illustrated in Figure 4.

Algorithm 1 NS Monge Iterations

Input: PSD matrix \mathbf{A}, \mathbf{B} , $\epsilon > 0$

$\mathbf{Y} \leftarrow \frac{\mathbf{B}}{(1+\epsilon)\|\mathbf{B}\|}$, $\mathbf{Z} \leftarrow \frac{\mathbf{A}}{(1+\epsilon)\|\mathbf{A}\|}$

while not converged **do**

$\mathbf{T} \leftarrow (3\mathbf{I} - \mathbf{Z}\mathbf{Y})/2$

$\mathbf{Y} \leftarrow \mathbf{Y}\mathbf{T}$

$\mathbf{Z} \leftarrow \mathbf{T}\mathbf{Z}$

end while

$\mathbf{Y} \leftarrow \sqrt{\frac{\|\mathbf{B}\|}{\|\mathbf{A}\|}}\mathbf{Y}$, $\mathbf{Z} \leftarrow \sqrt{\frac{\|\mathbf{A}\|}{\|\mathbf{B}\|}}\mathbf{Z}$

Output: $\mathbf{Y} = \mathbf{T}^{\mathbf{AB}}$, $\mathbf{Z} = \mathbf{T}^{\mathbf{BA}}$

Newton-Schulz iterations are quadratically convergent under the condition $\|\text{Id} - (\frac{\mathbf{A}}{\mathbf{0}} \frac{\mathbf{0}}{\mathbf{B}})^2\| < 1$, as shown in [30, Theorem 5.8]. To meet this condition, it is sufficient to rescale \mathbf{A} and \mathbf{B} so that their norms equal $(1 + \epsilon)^{-1}$ for some $\epsilon > 0$, as in the first step of Algorithm 1 (which can be skipped if $\|\mathbf{A}\| < 1$ (resp. $\|\mathbf{B}\| < 1$)). Finally, the output of the iterations are scaled back, using the homogeneity (resp. inverse homogeneity) of eq. (5) w.r.t. \mathbf{A} (resp. \mathbf{B}).

A rough theoretical analysis shows that both Newton-Schulz and SVD have a $O(d^3)$ complexity in the dimension. Figure 4 compares the running times of Newton-Schulz iterations and SVD on CPU or GPU used to compute both $\mathbf{A}^{\frac{1}{2}}$ and $\mathbf{A}^{-\frac{1}{2}}$. We simulate a batch of positive definite matrices \mathbf{A} following the Wishart distribution $W(\text{Id}_d, d)$ to which we add 0.1Id to avoid numerical issues when computing inverse square roots. We display the average run-time of 50 different trials along with its \pm std interval. Notice the different magnitudes between CPUs and GPUs. As a termination criterion, we first run EVD to obtain $\mathbf{A}_{evd}^{\frac{1}{2}}$ and $\mathbf{A}_{evd}^{-\frac{1}{2}}$ and stop the Newton-Schulz algorithm when its n -th running estimate $\mathbf{A}_n^{\frac{1}{2}}$ verifies: $\|\mathbf{A}_n^{\frac{1}{2}} - \mathbf{A}_{evd}^{\frac{1}{2}}\|_1 \leq 10^{-4}$. Notice the different order of magnitude between CPUs and GPUs. Moreover, the computational advantage of Newton-Schulz on GPUs can be further increased when computing multiple square roots in parallel.

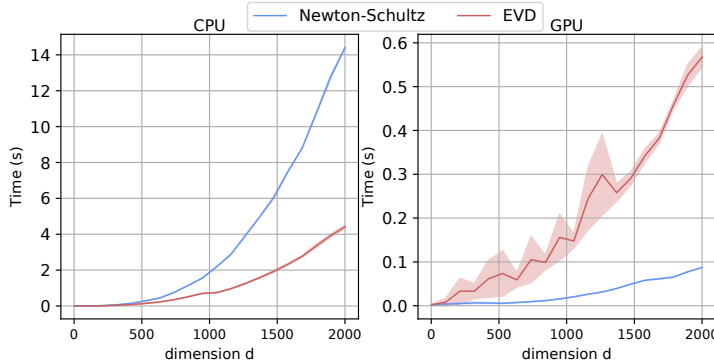


Figure 4: Average run-time of Newton-Schulz and EVD to compute on CPUs and GPUs.

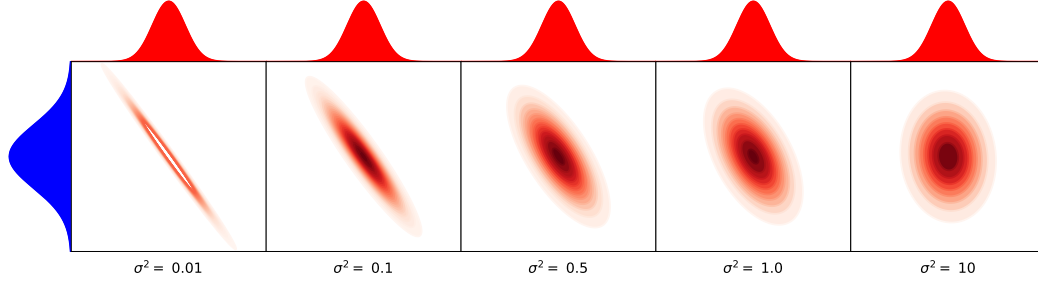


Figure 5: Effect of regularization on transportation plans. When σ goes to 0 (left), the transportation plan concentrates on the graph of the linear Monge map. When σ goes to infinity (right), the transportation plan converges to the independent coupling.

5.2 Effects of regularization strength.

We provide numerical experiments to illustrate the behaviour of transportation plans and corresponding distances as σ goes to 0 or to infinity. As can be seen from eq. (14), when $\sigma \rightarrow 0$ we recover the Wasserstein-Bures distance (3), and the optimal transportation plan converges to the Monge map (5). When on the contrary $\sigma \rightarrow \infty$, Sinkhorn divergences $\mathfrak{S}_\varepsilon(\alpha, \beta) \stackrel{\text{def}}{=} \text{OT}_\varepsilon(\alpha, \beta) - \frac{1}{2}(\text{OT}_\varepsilon(\alpha, \alpha) + \text{OT}_\varepsilon(\beta, \beta))$ convergence to MMD with a $-c$ kernel (where c is the optimal transport ground cost) [27]. With a $-\ell_2$ kernel, MMD is degenerate and equals 0 for centered measures.

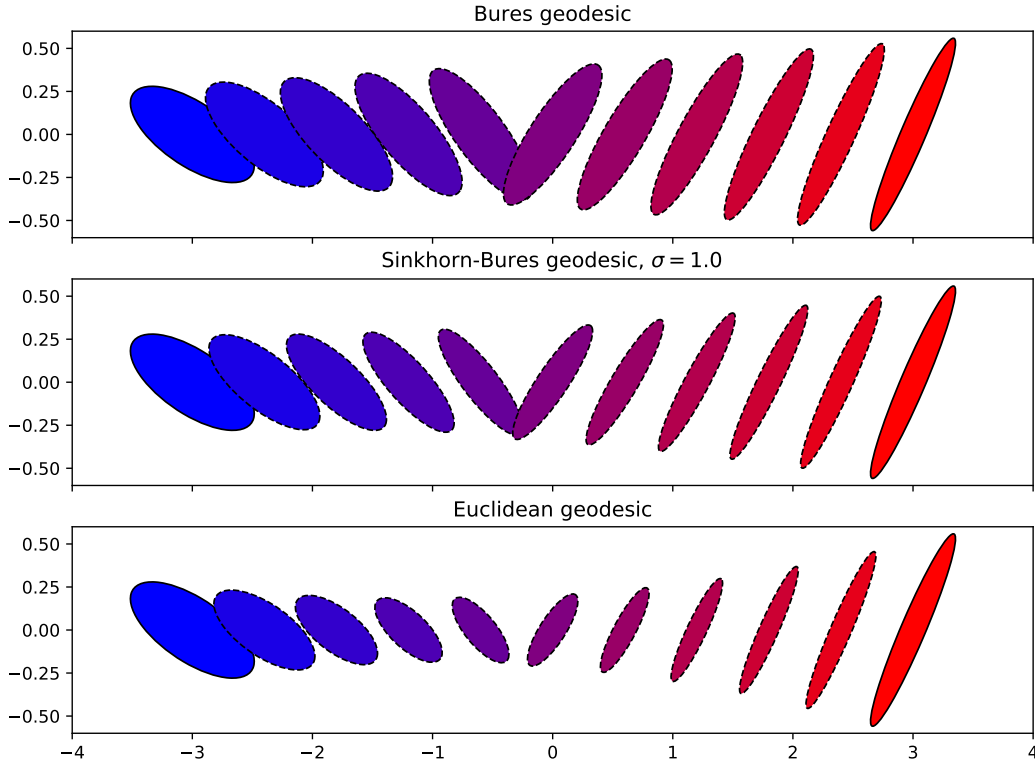


Figure 7: Bures, Sinkhorn-Bures, and Euclidean geodesics. Sinkhorn-Bures trajectories converge to Bures geodesics as σ goes to 0, and to Euclidean geodesics as σ goes to infinity.

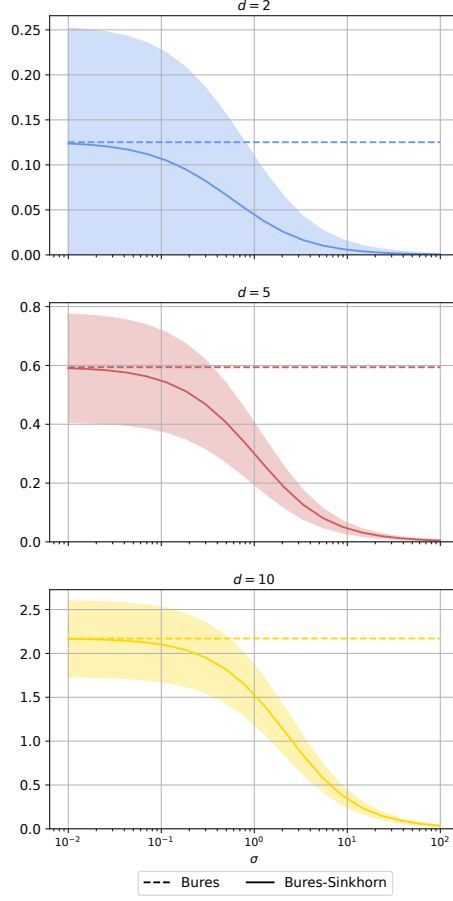


Figure 6: Numerical convergence of $\mathfrak{B}_\sigma^2(\mathbf{A}, \mathbf{B}) - \frac{1}{2}(\mathfrak{B}_\sigma^2(\mathbf{A}, \mathbf{A}) + \mathfrak{B}_\sigma^2(\mathbf{B}, \mathbf{B}))$ to $\mathfrak{B}^2(\mathbf{A}, \mathbf{B})$ as σ goes to 0 and to 0 as σ goes to infinity.

5.3 Proofs of technical results

We provide in this appendix the proofs of the results in the paper, as well as some technical lemmas used in solving Sinkhorn's equations in closed form.

Proof of Lemma 1.

Proof. Let $d\bar{\alpha}(x) = d\alpha(x + \mathbf{a})$ (resp. $d\bar{\beta}(y) = d\beta(y + \mathbf{b})$), $d\bar{\pi}(x, y) = d\pi(x + \mathbf{a}, y + \mathbf{b})$, such that $\bar{\alpha}, \bar{\beta}$ and $\bar{\pi}$ are centered. Then, $\forall \pi \in \Pi(\alpha, \beta)$,

- (i) $\bar{\pi} \in \Pi(\bar{\alpha}, \bar{\beta})$,
- (ii) $\text{KL}(\pi \| \alpha \otimes \beta) = \text{KL}(\bar{\pi} \| \bar{\alpha} \otimes \bar{\beta})$
- (iii) $\int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\bar{\pi}(x, y) = \int_{\mathbb{R}^d \times \mathbb{R}^d} \|(x - \mathbf{a}) - (y - \mathbf{b})\|^2 d\pi(x, y) = \|\mathbf{a} - \mathbf{b}\|^2 + \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y)$

Plugging (i)-(iii) into (7), we get $\text{OT}_\sigma(\alpha, \beta) = \text{OT}_\sigma(\bar{\alpha}, \bar{\beta}) + \|\mathbf{a} - \mathbf{b}\|^2$. □

Proof of Proposition 1.

Proof. The exponent inside the integral can be written as:

$$\begin{aligned} e^{-\frac{\|x-y\|^2}{2\sigma^2} + h(y)} d\alpha(y) &\propto e^{-\frac{\|x-y\|^2}{2\sigma^2} - \frac{1}{2}(y^\top \mathbf{X}y - y^\top \mathbf{A}^{-1}y)} dy \\ &\propto e^{-\frac{1}{2}(y^\top (\frac{\text{Id}}{\sigma^2} + \mathbf{X} + \mathbf{A}^{-1})y) + \frac{x^\top y}{\sigma^2}} dy \end{aligned}$$

which is integrable if and only if $\mathbf{X} + \mathbf{A}^{-1} + \frac{1}{\sigma^2} \text{Id} \succ 0$. Moreover, up to a multiplicative factor, the exponentiated Sinkhorn transform is equivalent to a Gaussian convolution of an exponentiated quadratic form. Lemma 4 applies:

$$\begin{aligned} e^{-T_\alpha(h)} &= \int_{\mathbb{R}^d} e^{-\frac{\|x-y\|^2}{2\sigma^2} + f(y)} d\alpha(y) \\ &\propto \int_{\mathbb{R}^d} e^{-\frac{\|x-y\|^2}{2\sigma^2} + \mathcal{Q}(\mathbf{X})(y) + \mathcal{Q}(\mathbf{A}^{-1})(y)} dy \\ &\propto \exp\left(\mathcal{Q}\left(\frac{\text{Id}}{\sigma^2}\right)\right) \star \exp\left(\mathcal{Q}(\mathbf{X}) + \mathcal{Q}(\mathbf{A}^{-1})\right) \\ &\propto \exp\left(\mathcal{Q}\left(\frac{\text{Id}}{\sigma^2}\right)\right) \star \exp\left(\mathcal{Q}(\mathbf{X} + \mathbf{A}^{-1})\right) \\ &\propto \exp\left(\mathcal{Q}((\text{Id} + \sigma^2\mathbf{X} + \sigma^2\mathbf{A}^{-1})^{-1}(\mathbf{X} + \mathbf{A}^{-1}))\right) \\ &\propto \exp\left(\mathcal{Q}\left(\frac{1}{\sigma^2}\mathbf{X}'^{-1}(\mathbf{X}' - \text{Id})\right)\right) \\ &\propto \exp\left(\mathcal{Q}\left(\frac{1}{\sigma^2}(\text{Id} - \mathbf{X}'^{-1})\right)\right). \end{aligned}$$

Therefore $T_\alpha(h)$ is up to an additive constant given by $\mathcal{Q}(\frac{1}{\sigma^2}(\mathbf{X}'^{-1} - \text{Id}))$.

Finally, since \mathbf{B} and \mathbf{X}' are positive definite, the positivity condition of \mathbf{Y}' holds and T_β can be applied again to get $T_\beta(T_\alpha(h))$. \square

Proof of Proposition 2.

Proof. Let $\mathbf{U}_0 = \mathbf{V}_0 = 0$. Applying Proposition 1 to the initial pair of potentials $\mathcal{Q}(\mathbf{U}_0), \mathcal{Q}(\mathbf{V}_0)$ leads to the sequence of quadratic Sinkhorn potentials $\frac{f_n}{2\sigma^2} = \mathcal{Q}(\mathbf{U}_n)$ and $\frac{g_n}{2\sigma^2} = \mathcal{Q}(\mathbf{V}_n)$ where:

$$\begin{aligned} \mathbf{V}_{n+1} &= \frac{1}{\sigma^2}((\sigma^2\mathbf{U}_n + \sigma^2\mathbf{A}^{-1} + \text{Id})^{-1} - \text{Id}) \\ \mathbf{U}_{n+1} &= \frac{1}{\sigma^2}((\sigma^2\mathbf{V}_{n+1} + \sigma^2\mathbf{B}^{-1} + \text{Id})^{-1} - \text{Id}). \end{aligned}$$

The change of variable:

$$\begin{aligned} \mathbf{F}_n &= \sigma^2\mathbf{U}_n + \sigma^2\mathbf{A}^{-1} + \text{Id} \\ \mathbf{G}_n &= \sigma^2\mathbf{V}_n + \sigma^2\mathbf{B}^{-1} + \text{Id} \end{aligned}$$

leads to (17).

We turn to show that this algorithm converges. First, note that since $\mathbf{F}_0, \mathbf{G}_0 \in \mathcal{S}_{++}^d$, a straightforward induction shows that $\forall n \geq 0, \mathbf{F}_n, \mathbf{G}_n \in \mathcal{S}_{++}^d$. Next, let us write the decoupled iteration on \mathbf{F} :

$$\mathbf{F} \leftarrow \sigma^2\mathbf{A}^{-1} + (\sigma^2\mathbf{B}^{-1} + \mathbf{F}^{-1})^{-1} \quad (42)$$

Let $\forall \mathbf{X} \in \mathcal{S}_{++}^d, \phi(\mathbf{X}) \stackrel{\text{def}}{=} \sigma^2\mathbf{A}^{-1} + (\sigma^2\mathbf{B}^{-1} + \mathbf{X}^{-1})^{-1} \in \mathcal{S}_{++}^d$. The first differential of ϕ admits the following expression:

$$\forall \mathbf{X} \in \mathcal{S}_{++}^d, \forall \mathbf{H} \in \mathbb{R}^{d \times d}, D\phi(\mathbf{X})[\mathbf{H}] = (\text{Id} + \sigma^2\mathbf{X}\mathbf{B}^{-1})^{-1}\mathbf{H}(\sigma^2\mathbf{B}^{-1}\mathbf{X} + \text{Id})^{-1}. \quad (43)$$

Hence, $\|D\phi(\mathbf{X})[\mathbf{H}]\|_{\text{op}} \leq \|(\text{Id} + \sigma^2\mathbf{X}\mathbf{B}^{-1})^{-1}\|_{\text{op}}^2 \|\mathbf{H}\|_{\text{op}}$. Plugging $\mathbf{H} = \text{Id}$, we get that $\|D\phi(\mathbf{X})\|_{\text{op}} = \|(\text{Id} + \sigma^2\mathbf{X}\mathbf{B}^{-1})^{-1}\|_{\text{op}}^2$. Finally, by matrix similarity

$$\|(\text{Id} + \sigma^2\mathbf{X}\mathbf{B}^{-1})^{-1}\|_{\text{op}} = \|(\text{Id} + \sigma^2\mathbf{B}^{-\frac{1}{2}}\mathbf{X}\mathbf{B}^{-\frac{1}{2}})^{-1}\|_{\text{op}} < 1,$$

which implies that $\|D\phi(\mathbf{X})\|_{\text{op}} < 1$ for $\mathbf{X} \in \mathcal{S}_{++}^d$ and $\sigma^2 > 0$. The same arguments hold for the iterates $(\mathbf{G}_n)_{n \geq 0}$.

From (42) and using Weyl's inequality, we can bound the smallest eigenvalue of \mathbf{F}_n from under: $\forall n, \lambda_d(\mathbf{F}_n) \geq \frac{\sigma^2}{\lambda_1(\mathbf{A})}$ (where $\lambda_d(\mathbf{F})$ is the smallest eigenvalue of \mathbf{F} and $\lambda_1(\mathbf{A})$ is the biggest eigenvalue of \mathbf{A}). Hence, the iterates live in $\mathcal{A} \stackrel{\text{def}}{=} \mathcal{S}_{++}^d \cap \{\mathbf{X} : \lambda_d(\mathbf{X}) \geq \frac{\sigma^2}{\lambda_1(\mathbf{A})}\}$. Finally, for all $\mathbf{X} \in \mathcal{A}$,

$$\begin{aligned} \|(\text{Id} + \sigma^2 \mathbf{B}^{-\frac{1}{2}} \mathbf{X} \mathbf{B}^{-\frac{1}{2}})^{-1}\|_{\text{op}} &= \frac{1}{\lambda_d(\text{Id} + \sigma^2 \mathbf{B}^{-1/2} \mathbf{X} \mathbf{B}^{-1/2})} \\ &= \frac{1}{1 + \sigma^2 \lambda_d(\mathbf{B}^{-1/2} \mathbf{X} \mathbf{B}^{-1/2})} \\ &\leq \frac{1}{1 + \sigma^2 \lambda_d(\mathbf{B}^{-1}) \lambda_d(\mathbf{X})} \\ &\leq \frac{1}{1 + \frac{\sigma^4}{\lambda_1(\mathbf{B}) \lambda_1(\mathbf{A})}} \end{aligned}$$

Which proves the uniform bound ■

□

Proof of Proposition 3.

Proof. Combining the two equations in (20) yields

$$\begin{aligned} \mathbf{G} &= \sigma^2 \mathbf{B}^{-1} + (\mathbf{G}^{-1} + \sigma^2 \mathbf{A}^{-1})^{-1} \\ \Leftrightarrow \mathbf{G} \mathbf{A}^{-1} &= \sigma^2 \mathbf{B}^{-1} \mathbf{A}^{-1} + (\mathbf{A} \mathbf{G}^{-1} + \sigma^2 \text{Id})^{-1} \\ \Leftrightarrow \mathbf{C}^{-1} &= \sigma^2 (\mathbf{A} \mathbf{B})^{-1} + (\mathbf{C} + \sigma^2 \text{Id})^{-1} \\ \Leftrightarrow \mathbf{C}^{-1} (\mathbf{C} + \sigma^2 \text{Id}) &= \sigma^2 (\mathbf{A} \mathbf{B})^{-1} (\mathbf{C} + \sigma^2 \text{Id}) + \text{Id} \\ \Leftrightarrow \text{Id} + \sigma^2 \mathbf{C}^{-1} &= \sigma^2 (\mathbf{A} \mathbf{B})^{-1} (\mathbf{C} + \sigma^2 \text{Id}) + \text{Id} \\ \Leftrightarrow \mathbf{C} + \sigma^2 \text{Id} &= \sigma^2 (\mathbf{A} \mathbf{B})^{-1} (\mathbf{C} + \sigma^2 \text{Id}) \mathbf{C} + \mathbf{C} \\ \Leftrightarrow \mathbf{C}^2 + \sigma^2 \mathbf{C} - \mathbf{A} \mathbf{B} &= 0. \end{aligned} \tag{44}$$

Given that \mathbf{A} and \mathbf{G}^{-1} are positive, their product $\mathbf{C} = \mathbf{A} \mathbf{G}^{-1}$ can be written: $\mathbf{A} \mathbf{G}^{-1} = \mathbf{A}^{\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{G}^{-1} \mathbf{A}^{\frac{1}{2}}) \mathbf{A}^{-\frac{1}{2}}$, thus $\mathbf{A} \mathbf{G}^{-1}$ is similar to the positive matrix $\mathbf{A}^{\frac{1}{2}} \mathbf{G}^{-1} \mathbf{A}^{\frac{1}{2}}$. Therefore, one can write an eigenvalue decomposition of $\mathbf{C} = \mathbf{P} \Sigma \mathbf{P}^{-1}$ with a positive diagonal matrix Σ . Substituting in (21), it follows that \mathbf{C} and $\mathbf{A} \mathbf{B}$ share the same eigenvectors with modified eigenvalues. Thus, it is sufficient to find the real roots of the polynomial $x \mapsto x^2 + \sigma^2 x - ab$ with $a, b \in \mathbb{R}_{++}$ which are given by: $x_1 = -\frac{\sigma^2}{2} - \sqrt{ab + \frac{\sigma^4}{4}}$ and $x_2 = -\frac{\sigma^2}{2} + \sqrt{ab + \frac{\sigma^4}{4}}$. Since \mathbf{C} is the product of the positive definite matrices \mathbf{G}^{-1} and \mathbf{A} , its eigenvalues are all positive. Discarding the negative root, the closed form follows immediately.

Indeed, by direct calculation, computing the square of the solution \mathbf{C} leads to the equation (21):

$$\begin{aligned} \mathbf{C}^2 &= \mathbf{A} \mathbf{B} + \frac{\sigma^4}{2} \text{Id} - \sigma^2 \left(\mathbf{A} \mathbf{B} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \\ &= \mathbf{A} \mathbf{B} - \sigma^2 \mathbf{C}. \end{aligned}$$

The second equality is obtained by observing that

$$(\mathbf{A}^{\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}})^2 = \mathbf{A}^{\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id}) \mathbf{A}^{-\frac{1}{2}} = \mathbf{A} \mathbf{B} + \frac{\sigma^4}{4} \text{Id},$$

i.e. that

$$\left(\mathbf{A} \mathbf{B} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} = \mathbf{A}^{\frac{1}{2}} (\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}}.$$

□

Proof of Lemma 2

Proof. It follows from elementary properties of Gaussian measures that the first and second marginals of π are respectively α and β . Hence,

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\pi(x, y) = \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x\|^2 d\pi(x, y) + \int_{\mathbb{R}^d \times \mathbb{R}^d} \|y\|^2 d\pi(x, y) - 2 \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle x, y \rangle d\pi(x, y) \quad (45)$$

$$= \int_{\mathbb{R}^d} \|x\|^2 d\alpha(x) + \int_{\mathbb{R}^d} \|y\|^2 d\beta(y) - 2 \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle x, y \rangle d\pi(x, y) \quad (46)$$

$$= \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - 2\text{Tr}(\mathbf{C}). \quad (47)$$

Next, using the closed form expression of the Kullback-Leibler divergence between Gaussian measures,

$$\text{KL}(\pi \| \alpha \otimes \beta) = \frac{1}{2} \left(\text{Tr} \left[\begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \right] - 2n + \log \det \begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{pmatrix} - \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \right) \quad (48)$$

$$= \frac{1}{2} (\log \det \mathbf{A} + \log \det \mathbf{B} - \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}). \quad (49)$$

□

Optimal transport plan and OT_σ

$$\begin{aligned} \frac{d\pi}{dx dy}(x, y) &= \exp \left(\frac{f(x) + g(y) - \|x - y\|^2}{2\sigma^2} \right) \frac{d\alpha}{dx}(x) \frac{d\beta}{dy}(y) \\ &\propto \exp \left(\mathcal{Q}(\mathbf{A}^{-1})(x) + \frac{f(x) + g(y) - \|x - y\|^2}{2\sigma^2} + \mathcal{Q}(\mathbf{B}^{-1})(y) \right) \\ &\propto \exp \left(\mathcal{Q}(\mathbf{U} + \mathbf{A}^{-1})(x) + \mathcal{Q}(\mathbf{V} + \mathbf{B}^{-1})(y) + \mathcal{Q} \left(\begin{pmatrix} \frac{\text{Id}}{\sigma^2} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\text{Id}}{\sigma^2} \end{pmatrix} \right)(x, y) \right) \\ &= \exp \left(\mathcal{Q} \begin{pmatrix} \mathbf{U} + \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} + \mathbf{B}^{-1} \end{pmatrix} (x, y) + \mathcal{Q} \left(\begin{pmatrix} \frac{\text{Id}}{\sigma^2} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\text{Id}}{\sigma^2} \end{pmatrix} \right)(x, y) \right) \\ &= \exp \left(\mathcal{Q} \begin{pmatrix} \frac{\text{Id}}{\sigma^2} + \mathbf{U} + \mathbf{A}^{-1} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\text{Id}}{\sigma^2} + \mathbf{V} + \mathbf{B}^{-1} \end{pmatrix} (x, y) \right) \\ &= \exp \left(\mathcal{Q} \begin{pmatrix} \frac{\mathbf{F}}{\sigma^2} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\mathbf{G}}{\sigma^2} \end{pmatrix} (x, y) \right) \\ &= \exp(\mathcal{Q}(\Gamma)(x, y)) \end{aligned}$$

with $\Gamma \stackrel{\text{def}}{=} \begin{pmatrix} \frac{\mathbf{F}}{\sigma^2} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\mathbf{G}}{\sigma^2} \end{pmatrix}$. Moreover, since $\frac{\mathbf{G}}{2\sigma^2} \succ 0$, and its Schur complement satisfies $\frac{\mathbf{F}}{\sigma^2} - \frac{1}{\sigma^2} \mathbf{G}^{-1} = \mathbf{A}^{-1} \succ 0$, we have that $\Gamma \succ 0$. Therefore π is a Gaussian $\mathcal{N}(\mathbf{H})$ with the covariance matrix given by the block inverse formula:

$$\mathbf{H} = \Gamma^{-1} \quad (50)$$

$$= \sigma^2 \begin{pmatrix} (\mathbf{F} - \mathbf{G}^{-1})^{-1} & (\mathbf{G}\mathbf{F} - \text{Id})^{-1} \\ (\mathbf{F}\mathbf{G} - \text{Id})^{-1} & (\mathbf{G} - \mathbf{F}^{-1})^{-1} \end{pmatrix} \quad (51)$$

$$= \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}, \quad (52)$$

where we used the optimality equations (20) and the definition of $\mathbf{C} = \mathbf{A}\mathbf{G}^{-1}$.

We can now conclude the proof of Theorem 1 by computing $\text{OT}_\sigma(\alpha, \beta)$ using Lemma 2. Let $\mathbf{R} = \mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}}$. Using the closed form expression of \mathbf{C} in (22), it first holds that

$$\mathbf{Z} \stackrel{\text{def}}{=} \mathbf{A}^{-\frac{1}{2}} \mathbf{C} \mathbf{A}^{\frac{1}{2}} = (\mathbf{R} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id}. \quad (53)$$

Moreover, since $\mathbf{R} = \mathbf{R}^\top$, it holds that $\mathbf{Z} = \mathbf{Z}^\top$. Hence,

$$\begin{aligned}
\det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} &= \det(\mathbf{A}) \det(\mathbf{B} - \mathbf{C}^\top \mathbf{A}^{-1} \mathbf{C}) \\
&= \det(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} - \mathbf{A}^{\frac{1}{2}} \mathbf{C}^\top \mathbf{A}^{-1} \mathbf{C} \mathbf{A}^{\frac{1}{2}}) \\
&= \det(\mathbf{R} - \mathbf{Z}^\top \mathbf{Z}) \\
&= \det(\mathbf{R} - \mathbf{Z}^2) \\
&= \det(\sigma^2(\mathbf{R} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} - \frac{\sigma^4}{2} \text{Id}) \\
&= (\frac{\sigma^2}{2})^d \det((4\mathbf{R} + \sigma^4 \text{Id})^{\frac{1}{2}} - \sigma^2 \text{Id}).
\end{aligned} \tag{54}$$

Since the matrices inside the determinant commute, we can use the identity $\mathbf{P} - \mathbf{Q} = (\mathbf{P}^2 - \mathbf{Q}^2)(\mathbf{P} + \mathbf{Q})^{-1}$ to get rid of the negative sign. Equation (54) then becomes:

$$\begin{aligned}
(\frac{\sigma^2}{2})^d \det((4\mathbf{R} + \sigma^4 \text{Id})^{\frac{1}{2}} - \sigma^2 \text{Id}) &= (\frac{\sigma^2}{2})^d \det(4\mathbf{R}) \det\left(\left((4\mathbf{R} + \sigma^4 \text{Id})^{\frac{1}{2}} + \sigma^2 \text{Id}\right)^{-1}\right) \\
&= (2\sigma^2)^d \det(\mathbf{A}\mathbf{B}) \det\left(\left((4\mathbf{R} + \sigma^4 \text{Id})^{\frac{1}{2}} + \sigma^2 \text{Id}\right)^{-1}\right).
\end{aligned}$$

Plugging this expression in (25), the determinant of \mathbf{A} and \mathbf{B} cancel out and we finally get:

$$\begin{aligned}
\mathfrak{B}_\sigma^2(\mathbf{A}, \mathbf{B}) &= \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - \text{Tr}(4\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} + \sigma^4 \text{Id})^{\frac{1}{2}} + d\sigma^2 - \\
&\quad \sigma^2 d \log(2\sigma^2) + \sigma^2 \log \det\left(\left(4\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} + \sigma^4 \text{Id}\right)^{\frac{1}{2}} + \sigma^2 \text{Id}\right).
\end{aligned}$$

Proof of Proposition 4

Proof. Using Lemma 2, eq. (7) becomes

$$\mathfrak{B}_\sigma^2(\mathbf{A}, \mathbf{B}) = \min_{\mathbf{C}: \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \geq 0} \left\{ \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) - 2\text{Tr}(\mathbf{C}) + \sigma^2(\log \det \mathbf{A} + \log \det \mathbf{B} - \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}) \right\},$$

which gives eq. (26). Let us now prove eq. (27). A necessary and sufficient condition for $\begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \geq 0$

is that there exists a contraction \mathbf{K} (i.e. $\mathbf{K} \in \mathbb{R}^d : \|\mathbf{K}\|_{\text{op}} \leq 1$) such that $\mathbf{C} = \mathbf{A}^{\frac{1}{2}} \mathbf{K} \mathbf{B}^{\frac{1}{2}}$ [4, Ch. 1].¹ With this parameterization, we have (using Schur complements) that

$$\begin{aligned}
\det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} &= \det \mathbf{B} \det(\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}^\top) \\
&= \det \mathbf{B} \det \mathbf{A} \det(\text{Id} - \mathbf{K} \mathbf{K}^\top)
\end{aligned}$$

Hence, injecting this in Equation (26), we have the following equivalent problem:

$$\mathfrak{B}_\sigma^2(\mathbf{A}, \mathbf{B}) = \min_{\mathbf{K} \in \mathbb{R}^{d \times d} : \|\mathbf{K}\|_{\text{op}} \leq 1} \text{Tr} \mathbf{A} + \text{Tr} \mathbf{B} - 2\text{Tr} \mathbf{A}^{\frac{1}{2}} \mathbf{K} \mathbf{B}^{\frac{1}{2}} - \sigma^2 \ln \det(\text{Id} - \mathbf{K} \mathbf{K}^\top) \tag{55}$$

Let's prove that both problems are convex.

- (26): The set $\{\mathbf{C} : \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix} \geq 0\}$ is convex, since $\begin{pmatrix} \mathbf{A} & \mathbf{C}_1 \\ \mathbf{C}_1^\top & \mathbf{B} \end{pmatrix} \geq 0$ and $\begin{pmatrix} \mathbf{A} & \mathbf{C}_2 \\ \mathbf{C}_2^\top & \mathbf{B} \end{pmatrix} \geq 0$ implies that $\begin{pmatrix} \mathbf{A} & (1-\theta)\mathbf{C}_1 + \theta\mathbf{C}_2 \\ ((1-\theta)\mathbf{C}_1^\top + \theta\mathbf{C}_2^\top) & \mathbf{B} \end{pmatrix} = (1-\theta) \begin{pmatrix} \mathbf{A} & \mathbf{C}_1 \\ \mathbf{C}_1^\top & \mathbf{B} \end{pmatrix} + \theta \begin{pmatrix} \mathbf{A} & \mathbf{C}_2 \\ \mathbf{C}_2^\top & \mathbf{B} \end{pmatrix} \geq 0$. Following the same decomposition, the concavity of the log det function implies that $\mathbf{C} \rightarrow \log \det \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{pmatrix}$ is concave, and hence that the objective function of (26) is convex.
- (27): The ball $\mathcal{B}_{\text{op}} \stackrel{\text{def}}{=} \{\mathbf{K} \in \mathbb{R}^{d \times d} : \|\mathbf{K}\|_{\text{op}} \leq 1\}$ is obviously convex. Hence, there remains to prove that $f(\mathbf{K}) : \mathbf{K} \in \mathcal{B}_{\text{op}} \rightarrow \log \det(\text{Id} - \mathbf{K} \mathbf{K}^\top)$ is concave. Indeed, it holds that $f(\mathbf{K}) = \log \det \begin{pmatrix} \text{Id} & \mathbf{K} \\ \mathbf{K}^\top & \text{Id} \end{pmatrix}$. Hence, $\forall \mathbf{K}, \mathbf{H} \in \mathcal{B}_{\text{op}}, \forall t \in [0, 1]$,

$$\begin{aligned}
f((1-t)\mathbf{K} + t\mathbf{H}) &= \log \det \left\{ (1-t) \begin{pmatrix} \text{Id} & \mathbf{K} \\ \mathbf{K}^\top & \text{Id} \end{pmatrix} + t \begin{pmatrix} \text{Id} & \mathbf{H} \\ \mathbf{H}^\top & \text{Id} \end{pmatrix} \right\} \\
&\geq (1-t) \log \det \begin{pmatrix} \text{Id} & \mathbf{K} \\ \mathbf{K}^\top & \text{Id} \end{pmatrix} + t \log \det \begin{pmatrix} \text{Id} & \mathbf{H} \\ \mathbf{H}^\top & \text{Id} \end{pmatrix} \\
&= (1-t)f(\mathbf{K}) + tf(\mathbf{H}),
\end{aligned}$$

where the second line follows from the concavity of log det.

¹Another immediate NSC is $\mathbf{A} \geq \mathbf{C} \mathbf{B}^{-1} \mathbf{C}^\top$

□

Proof of Proposition 5

Proof. By Proposition 4, (26) is convex, hence strong duality holds. Ignoring the terms not depending on \mathbf{C} , problem (26) can be written using the redundant parameterization $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 & \mathbf{X}_2 \\ \mathbf{X}_3 & \mathbf{X}_4 \end{pmatrix}$:

$$\mathfrak{D}(\mathbf{A}, \mathbf{B}) \stackrel{\text{def}}{=} \min_{\substack{\mathbf{X} \succ 0 \\ \mathbf{X}_1 = \mathbf{A}, \mathbf{X}_4 = \mathbf{B}}} -\text{Tr}(\mathbf{X}_2) - \text{Tr}(\mathbf{X}_3) - \sigma^2 \log \det(\mathbf{X}) \quad (56)$$

$$= \min_{\substack{\mathbf{X} \succ 0 \\ \mathbf{X}_1 = \mathbf{A}, \mathbf{X}_4 = \mathbf{B}}} -\langle \mathbf{X}, \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \rangle - \sigma^2 \log \det(\mathbf{X}) \quad (57)$$

$$= \min_{\substack{\mathbf{X} \succ 0 \\ \mathbf{X}_1 = \mathbf{A}, \mathbf{X}_4 = \mathbf{B}}} \mathcal{F}(\mathbf{X}), \quad (58)$$

where the functional \mathcal{F} is convex. Moreover, its Legendre transform is given by:

$$\begin{aligned} \mathcal{F}^*(\mathbf{Y}) &= \max_{\mathbf{X} \succ 0} \langle \mathbf{X}, \mathbf{Y} + \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \rangle + \sigma^2 \log \det(\mathbf{X}) \\ &= (-\sigma^2 \log \det)^* \left(\mathbf{Y} + \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \right) \\ &= \sigma^2 (-\log \det)^* \left(\frac{1}{\sigma^2} \left(\mathbf{Y} + \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \right) \right) \\ &= -\sigma^2 \log \det \left(-\frac{1}{\sigma^2} \left(\mathbf{Y} + \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \right) \right) - 2\sigma^2 d \\ &= -\sigma^2 \log \det \left(- \left(\mathbf{Y} + \begin{pmatrix} 0 & \text{Id} \\ \text{Id} & 0 \end{pmatrix} \right) \right) - 2d(\sigma^2 - \sigma^2 \log(\sigma^2)). \end{aligned}$$

Let \mathcal{H} be the linear operator $\mathcal{H} : \mathbf{X} \mapsto (\mathbf{X}_1, \mathbf{X}_4)$. Its conjugate operator is defined on $\mathcal{S}_{++}^d \times \mathcal{S}_{++}^d$ and is given by $\mathcal{H}^*(\mathbf{F}, \mathbf{G}) = \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix}$. Therefore, Fenchel's duality theorem leads to:

$$\begin{aligned} \mathfrak{D}(\mathbf{A}, \mathbf{B}) &= \max_{\mathbf{F}, \mathbf{G} \succ 0} -\langle \mathbf{F}, \mathbf{A} \rangle - \langle \mathbf{G}, \mathbf{B} \rangle - \mathcal{F}^*(-\mathcal{H}^*(\mathbf{F}, \mathbf{G})) \\ &= \max_{\mathbf{F}, \mathbf{G} \succ 0} -\langle \mathbf{F}, \mathbf{A} \rangle - \langle \mathbf{G}, \mathbf{B} \rangle + \sigma^2 \log \det \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix} + 2d(\sigma^2 - \sigma^2 \log(\sigma^2)) \\ &= \max_{\mathbf{F}, \mathbf{G} \succ 0} -\langle \mathbf{F}, \mathbf{A} \rangle - \langle \mathbf{G}, \mathbf{B} \rangle + \sigma^2 \log \det(\mathbf{FG} - \text{Id}) + 2d(\sigma^2 - \sigma^2 \log(\sigma^2)) \end{aligned}$$

Where the last equality follows from the fact that Id and \mathbf{G} commute. Therefore, reinserting the discarded trace terms, the dual problem of (26) can be written as

$$\begin{aligned} \max_{\mathbf{F}, \mathbf{G} \succ 0} \Big\{ & -\langle \mathbf{F}, \mathbf{A} \rangle - \langle \mathbf{G}, \mathbf{B} \rangle + \sigma^2 \log \det(\mathbf{FG} - \text{Id}) \\ & + \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}) + \sigma^2 \log \det \mathbf{AB} + 2d\sigma^2(1 - \log \sigma^2) \Big\}. \end{aligned} \quad (59)$$

□

Proof of Proposition 6

Proof. (i) *Optimality:* Canceling out the gradients in eq. (28) leads to the following optimality conditions:

$$\begin{aligned} -\mathbf{A} + \sigma^2 \mathbf{G}(\mathbf{FG} - \text{Id})^{-1} &= 0 \\ -\mathbf{B} + \sigma^2 (\mathbf{FG} - \text{Id})^{-1} \mathbf{F} &= 0, \end{aligned} \quad (60)$$

i.e.

$$\begin{aligned} \mathbf{F} &= \sigma^2 \mathbf{A}^{-1} + \mathbf{G}^{-1} \\ \mathbf{G} &= \sigma^2 \mathbf{B}^{-1} + \mathbf{F}^{-1} \end{aligned} \quad (61)$$

Thus (\mathbf{F}, \mathbf{G}) is a solution of the Sinkhorn fixed point equation (20).

(ii) *Differentiability*: Using Danskin's theorem on problem (28) leads to the formula of the gradient as a function of the optimal dual pair (\mathbf{F}, \mathbf{G}) . Indeed, keeping in mind that $\nabla_{\mathbf{A}} \log \det(\mathbf{A}) = -\mathbf{A}^{-1}$ and using the change of variable of Proposition 2, we recover the dual potentials of Corollary 1:

$$\begin{aligned}\nabla \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B}) &= (\text{Id} - \mathbf{F}^* + \sigma^2 \mathbf{A}^{-1}, \text{Id} - \mathbf{G}^* + \sigma^2 \mathbf{B}^{-1}) \\ &= -\sigma^2(\mathbf{U}, \mathbf{V})\end{aligned}$$

Using Corollary 1, it holds that

$$\begin{aligned}\nabla_{\mathbf{A}} \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B}) &= -\sigma^2 \mathbf{U} \\ &= \text{Id} - \mathbf{B}(\mathbf{C} + \sigma^2 \text{Id})^{-1} \\ &= \text{Id} - \mathbf{B} \left((\mathbf{A}\mathbf{B} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \\ &= \text{Id} - \mathbf{B}^{\frac{1}{2}} \left((\mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}} \\ &= \text{Id} - \mathbf{B}^{\frac{1}{2}} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}},\end{aligned}$$

where $\mathbf{D} \stackrel{\text{def}}{=} \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id}$.

(iii) *Convexity*: Assume without loss of generality that \mathbf{B} is fixed and let $G : \mathbf{B} \mapsto \nabla_{\mathbf{A}} \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$. As long as $\sigma > 0$, G is differentiable as a composition of differentiable functions. Let's show that the Hessian of $\psi : \mathbf{A} \mapsto \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$ is a positive quadratic form. Take a direction $\mathbf{H} \in \mathcal{S}_+^d$. It holds:

$$\begin{aligned}\nabla_{\mathbf{A}}^2 \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})(\mathbf{H}, \mathbf{H}) &= \langle \mathbf{H}, \text{Jac}_G(\mathbf{A})(\mathbf{H}) \rangle \\ &= \text{Tr}(\mathbf{H} \text{Jac}_G(\mathbf{A})(\mathbf{H})).\end{aligned}$$

For the sake of clarity, let's write $G(\mathbf{A}) = \text{Id} - L(W(\phi(\mathbf{A})))$ with the following intermediary functions:

$$\begin{aligned}L : \mathbf{A} &\mapsto \mathbf{B}^{\frac{1}{2}} \mathbf{A} \mathbf{B}^{\frac{1}{2}} \\ Q : \mathbf{A} &\mapsto \mathbf{A}^{\frac{1}{2}} \\ \phi : \mathbf{A} &\mapsto Q(L(\mathbf{A}) + \frac{\sigma^4}{4} \text{Id}) \\ W : \mathbf{A} &\mapsto (\mathbf{A} + \frac{\sigma^2}{2} \text{Id})^{-1}.\end{aligned}$$

Moreover, their derivatives are given by:

$$\begin{aligned}\text{Jac}_L(\mathbf{A})(\mathbf{H}) &= \mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}} \\ \text{Jac}_W(\mathbf{A})(\mathbf{H}) &= -(\mathbf{A} + \frac{\sigma^2}{2} \text{Id})^{-1} \mathbf{H} (\mathbf{A} + \frac{\sigma^2}{2} \text{Id})^{-1} \\ \text{Jac}_Q(\mathbf{A})(\mathbf{H}) &= \mathbf{Z},\end{aligned}$$

where $\mathbf{Z} \in \mathcal{S}_+^d$ is the unique solution of the Sylvester equation: $\mathbf{Z} \mathbf{A}^{\frac{1}{2}} + \mathbf{A}^{\frac{1}{2}} \mathbf{Z} = \mathbf{H}$.

Using the chain rule:

$$\begin{aligned}\text{Jac}_G(\mathbf{A})(\mathbf{H}) &= -\text{Jac}_L(W(\phi(\mathbf{A}))) (\text{Jac}_W(\phi(\mathbf{A})) (\text{Jac}_\phi(\mathbf{A})(\mathbf{H}))) \\ &= -\mathbf{B}^{\frac{1}{2}} \text{Jac}_W(\phi(\mathbf{A})) (\text{Jac}_\phi(\mathbf{A})(\mathbf{H})) \mathbf{B}^{\frac{1}{2}} \\ &= \mathbf{B}^{\frac{1}{2}} \left(\phi(\mathbf{A}) + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \text{Jac}_\phi(\mathbf{A})(\mathbf{H}) \left(\phi(\mathbf{A}) + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}} \\ &= \mathbf{B}^{\frac{1}{2}} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \text{Jac}_\phi(\mathbf{A})(\mathbf{H}) \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}}.\end{aligned}$$

Again using the chain rule:

$$\begin{aligned}
\mathbf{Y} &\stackrel{\text{def}}{=} \text{Jac}_\phi(\mathbf{A})(\mathbf{H}) = \text{Jac}_Q(L(\mathbf{A}) + \frac{\sigma^4}{4} \text{Id})(\text{Jac}_L(\mathbf{A})(\mathbf{H})) \\
&= \text{Jac}_Q(L(\mathbf{A}) + \frac{\sigma^4}{4} \text{Id})(\mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}}) \\
&= \text{Jac}_Q(\mathbf{D})(\mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}}).
\end{aligned}$$

Therefore, $\mathbf{Y} \succ 0$ is the unique solution of the Sylvester equation:

$$\mathbf{Y} \mathbf{D}^{\frac{1}{2}} + \mathbf{D}^{\frac{1}{2}} \mathbf{Y} = \mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}}.$$

Combining everything:

$$\begin{aligned}
\nabla_{\mathbf{A}}^2 \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})(\mathbf{H}, \mathbf{H}) &= \langle \mathbf{H}, \text{Jac}_G(\mathbf{A})(\mathbf{H}) \rangle \\
&= \text{Tr}(\mathbf{H} \text{Jac}_G(\mathbf{A})(\mathbf{H})) \\
&= \text{Tr} \left(\mathbf{H} \mathbf{B}^{\frac{1}{2}} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{Y} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{B}^{\frac{1}{2}} \right) \\
&= \text{Tr} \left(\mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{Y} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \right).
\end{aligned}$$

Since \mathbf{H} and \mathbf{Y} are positive, the matrices $\mathbf{B}^{\frac{1}{2}} \mathbf{H} \mathbf{B}^{\frac{1}{2}}$ and $\left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{Y} \left(\mathbf{D}^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1}$ are positive semi-definite as well. Their product is similar to a positive semi-definite matrix, therefore the trace above is non-negative.

Given that \mathbf{A} and \mathbf{H} are arbitrary positive semi-definite matrices, it holds that

$$\nabla_{\mathbf{A}}^2 \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})(\mathbf{H}, \mathbf{H}) \geq 0$$

Therefore, $\mathbf{A} \mapsto \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$ is convex.

Counter-example of joint convexity: If \mathfrak{B}_{σ^2} were jointly convex, then $\delta \stackrel{\text{def}}{=} : \mathbf{A} \rightarrow \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{A})$ would be a convex function.

In the 1-dimensional case with $\sigma = 1$, one can see that this would be equivalent to $x \rightarrow \ln((x^2 + 1)^{\frac{1}{2}} + 1) - (x^2 + 1)^{\frac{1}{2}}$ being convex, whereas it is in fact strictly concave.

(iv) *Minimizer of $\phi_{\mathbf{B}}$* With fixed \mathbf{B} , cancelling the gradient of $\phi_{\mathbf{B}} \stackrel{\text{def}}{=} : \mathbf{A} \mapsto \mathfrak{B}_{\sigma^2}(\mathbf{A}, \mathbf{B})$ leads to $\mathbf{A} = \mathbf{B} - \sigma^2 \text{Id}$ which is well defined if and only if $\mathbf{B} \succeq \sigma^2 \text{Id}$. However, if $\mathbf{B} - \sigma^2 \text{Id}$ is not positive semi-definite, write the eigenvalue decomposition: $\mathbf{B} = \mathbf{P} \Sigma \mathbf{P}^\top$ and define $\mathbf{A}_0 \stackrel{\text{def}}{=} \mathbf{P}(\Sigma - \sigma^2 \text{Id})_+ \mathbf{P}^\top$ where the operator $x_+ = \max(x, 0)$ is applied element-wise. Then:

$$\begin{aligned}
\nabla_{\mathbf{A}} \phi_{\mathbf{B}}(\mathbf{A}_0) &= \text{Id} - \mathbf{P} \Sigma^{\frac{1}{2}} \mathbf{P}^\top \left((\mathbf{P}(\Sigma^2 - \sigma^2 \Sigma)_+ \mathbf{P}^\top + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \mathbf{P} \Sigma^{\frac{1}{2}} \mathbf{P}^\top \\
&= \text{Id} - \mathbf{P} \Sigma^{\frac{1}{2}} \left(((\Sigma^2 - \sigma^2 \Sigma)_+ + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2} \text{Id} \right)^{-1} \Sigma^{\frac{1}{2}} \mathbf{P}^\top \\
&= \text{Id} - \mathbf{P} \Sigma^{\frac{1}{2}} ((\Sigma - \sigma^2 \text{Id})_+ + \sigma^2 \text{Id})^{-1} \Sigma^{\frac{1}{2}} \mathbf{P}^\top \\
&= \mathbf{P}(\text{Id} - \Sigma^{\frac{1}{2}} ((\Sigma - \sigma^2 \text{Id})_+ + \sigma^2 \text{Id})^{-1} \Sigma^{\frac{1}{2}}) \mathbf{P}^\top \\
&= \frac{1}{\sigma^2} \mathbf{P}(\sigma^2 \text{Id} - \Sigma)_+ \mathbf{P}^\top
\end{aligned}$$

Thus, given that $(\Sigma - \sigma^2 \text{Id})_+(\sigma^2 \text{Id} - \Sigma)_+ = 0$, it holds, for any $\mathbf{H} \in \mathcal{S}_+^d$:

$$\begin{aligned}
\langle \mathbf{H} - \mathbf{A}_0, \nabla_{\mathbf{A}} \phi_{\mathbf{B}}(\mathbf{A}_0) \rangle &= \langle \mathbf{P}^\top \mathbf{H} \mathbf{P} - (\Sigma - \sigma^2 \text{Id})_+, (\sigma^2 \text{Id} - \Sigma)_+ \rangle \\
&= \langle \mathbf{P}^\top \mathbf{H} \mathbf{P}, (\sigma^2 \text{Id} - \Sigma)_+ \rangle \\
&= \text{Tr}(\mathbf{P}^\top \mathbf{H} \mathbf{P}(\sigma^2 \text{Id} - \Sigma)_+) \geq 0
\end{aligned}$$

Where the last inequality holds since both matrices are positive semi-definite. Given that $\phi_{\mathbf{B}}$ is convex, the first order optimality condition holds so $\phi_{\mathbf{B}}$ is minimized at \mathbf{A}_0 . \square

Proof of Theorem 2

Proof. This theorem is a generalization of [31, Thm 3] for multivariate Gaussians. First we are going to break it down using the centering lemma 1. For any probability measure μ , let $\bar{\mu}$ denote its centered transformation. The debiased barycenter problem is equivalent to:

$$\begin{aligned}
& \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k S_{\sigma}(\alpha_k, \beta) \\
&= \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \text{OT}_{\sigma}(\alpha_k, \beta) - \frac{1}{2}(\text{OT}_{\sigma}(\alpha_k, \alpha_k) + \text{OT}_{\sigma}(\beta, \beta)) \\
&= \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \|\mathbf{a}_k - \mathbb{E}_{\beta}(X)\|^2 + w_k \text{OT}_{\sigma}(\bar{\alpha}_k, \bar{\beta}) - \frac{1}{2}(w_k \text{OT}_{\sigma}(\bar{\alpha}_k, \bar{\alpha}_k) + \text{OT}_{\sigma}(\bar{\beta}, \bar{\beta})) \\
&= \min_{\substack{\mathbf{b} \in \mathbb{R}^d \\ \beta \in \mathcal{G}, \mathbb{E}_{\beta}(\mathbf{X})=0}} \sum_{k=1}^K w_k \|\mathbf{a}_k - \mathbf{b}\|^2 + w_k \text{OT}_{\sigma}(\bar{\alpha}_k, \beta) - \frac{1}{2}(w_k \text{OT}_{\sigma}(\bar{\alpha}_k, \bar{\alpha}_k) + \text{OT}_{\sigma}(\beta, \beta))
\end{aligned} \tag{62}$$

Therefore, since both arguments are independent, we can first minimize over \mathbf{b} to obtain $\mathbb{E}_{\beta}(\mathbf{X}) = \mathbf{b} = \sum_{k=1}^K w_k \mathbf{a}_k$. Without loss of generality, we assume from now on that $\mathbf{a}_k = 0$ for all k .

The rest of this proof is adapted from [31], Thm 3 to $d \geq 1$. Janati et al. [31] showed that S_{σ} is differentiable and convex (w.r.t. one measure at a time) on sub-Gaussian measures where the notion of differentiability is different from the usual Fréchet differentiability: a function $F : \mathcal{G} \rightarrow \mathbb{R}$ is differentiable at α if there exists $\nabla F(\alpha) \in \mathcal{C}(\mathbb{R}^d)$ such that for any displacement $t\delta\alpha$ with $t > 0$ and $\delta\alpha = \alpha_1 - \alpha_2$ with $\alpha_1, \alpha_2 \in \mathcal{G}$, and

$$F(\alpha + t\delta\alpha) = F(\alpha) + t\langle \delta\alpha, \nabla F(\alpha) \rangle + o(t) , \tag{63}$$

where $\langle \delta\alpha, \nabla F(\alpha) \rangle = \int_{\mathbb{R}^d} \nabla F(\alpha) d\delta\alpha$.

Moreover, F is convex if and only if for any $\alpha, \alpha' \in \mathcal{G}$:

$$F(\alpha) \geq F(\alpha') + \langle \alpha - \alpha', \nabla F(\alpha') \rangle , \tag{64}$$

Let (f_k, g_k) denote the potentials associated with $\text{OT}_{\sigma}(\alpha_k, \beta)$ and h_{β} the autocorrelation potential associated with $\text{OT}_{\sigma}(\beta, \beta)$. If β is sub-Gaussian, it holds: $\nabla_{\beta} S_{\sigma}(\alpha_k, \beta) = g_k - h$. Therefore, from (64) a probability measure β is the debiased barycenter if and only if for any direction $\mu \in \mathcal{G}$, the optimality condition holds:

$$\begin{aligned}
& \left\langle \sum_{k=1}^K w_k \nabla_{\beta} S_{\sigma}(\alpha_k, \beta), \mu - \beta \right\rangle \geq 0 \\
& \Leftrightarrow \sum_{k=1}^K w_k \langle g_k - h_{\beta}, \mu - \beta \rangle \geq 0
\end{aligned} \tag{65}$$

Moreover, the potentials $(f_k), (g_k)$ and h must verify the Sinkhorn optimality conditions (10) for all k and for all x β -a.s and y α -a.s:

$$\begin{cases} e^{\frac{f_k(x)}{2\sigma^2}} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2 + g_k(y)}{2\sigma^2}} d\beta(y) \right) = 1, & e^{\frac{g_k(x)}{2\sigma^2}} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2 + f_k(y)}{2\sigma^2}} d\alpha_k(y) \right) = 1. \\ e^{\frac{h(x)}{2\sigma^2}} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2 + h_{\beta}(y)}{2\sigma^2}} d\beta(y) \right) = 1. \end{cases} \tag{66}$$

We are going to show that for the Gaussian measure β given in the statement of the theorem is well-defined and verifies all optimality conditions (66). Indeed, assume that β is a Gaussian measure given by $\mathcal{N}(\mathbf{B})$ for some unknown $\mathbf{B} \in S_+^d$ (remember that β is necessarily centered, following the developments (62)). The Sinkhorn equations can therefore be written as a system on positive definite matrices:

$$\mathbf{F}_k = \sigma^2 \mathbf{A}_k^{-1} + \mathbf{G}_k^{-1}, \quad \mathbf{G}_k = \sigma^2 \mathbf{B} + \mathbf{F}_k^{-1}, \quad \mathbf{H} = \sigma^2 \mathbf{B} + \mathbf{H}^{-1}$$

where for all k :

$$\begin{aligned}\frac{f_k}{2\sigma^2} &= \mathcal{Q}\left(\frac{1}{\sigma^2}(\mathbf{G}_k^{-1} - \text{Id})\right) + f_k(0) \\ \frac{g_k}{2\sigma^2} &= \mathcal{Q}\left(\frac{1}{\sigma^2}(\mathbf{F}_k^{-1} - \text{Id})\right) + g_k(0) \\ \frac{h}{2\sigma^2} &= \mathcal{Q}\left(\frac{1}{\sigma^2}(\mathbf{H}^{-1} - \text{Id})\right) + h_\beta(0)\end{aligned}\tag{67}$$

Moreover, provided \mathbf{B} exists and is positive definite, the system (67) has a unique set of solutions $(\mathbf{F}_k)_k, (\mathbf{G}_k)_k, \mathbf{H}$ given by:

$$\mathbf{F}_k = \mathbf{B}\mathbf{C}_k^{-1}, \quad \mathbf{G}_k = \mathbf{C}_k^{-1}\mathbf{A}_k, \quad \mathbf{H} = \mathbf{B}^{-1}\mathbf{J}\tag{68}$$

where $\mathbf{C}_k = (\mathbf{A}_k\mathbf{B} + \frac{\sigma^4}{4}\text{Id})^{\frac{1}{2}} - \frac{\sigma^2}{2}\text{Id}$ and $\mathbf{J} = (\mathbf{B}^2 + \frac{\sigma^4}{4}\text{Id})^{\frac{1}{2}} + \frac{\sigma^2}{2}\text{Id}$. Therefore, the gradient in (65) can be written:

$$\begin{aligned}\sum_{k=1}^K w_k \langle g_k - h_\beta \rangle &= \mathcal{Q}\left(\frac{1}{\sigma^2}\left(\sum_{k=1}^K w_k \mathbf{F}_k^{-1} - \mathbf{H}^{-1}\right)\right) + \sum_{w=1}^K w_k g_k(0) - h_\beta(0) \\ &\propto \mathcal{Q}\left(\sum_{k=1}^K w_k \mathbf{C}_k \mathbf{B}^{-1} - \mathbf{J}^{-1} \mathbf{B}\right) + \sum_{w=1}^K w_k g_k(0) - h_\beta(0)\end{aligned}\tag{69}$$

and

$$\begin{aligned}&\sum_{k=1}^K w_k \mathbf{C}_k \mathbf{B}^{-1} - \mathbf{J}^{-1} \mathbf{B} \\ &= \sum_{k=1}^K w_k \mathbf{B}^{-\frac{1}{2}} \left(\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \mathbf{B}^{-\frac{1}{2}} - \mathbf{B}^{-1} \left(\mathbf{B}^2 + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \\ &= \sum_{k=1}^K w_k \mathbf{B}^{-\frac{1}{2}} \left(\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \mathbf{B}^{-\frac{1}{2}} - \mathbf{B}^{-\frac{1}{2}} \left(\mathbf{B}^2 + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \mathbf{B}^{-\frac{1}{2}} \\ &= \mathbf{B}^{-\frac{1}{2}} \left(\sum_{k=1}^K w_k \left(\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \left(\mathbf{B}^2 + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \right) \mathbf{B}^{-\frac{1}{2}}\end{aligned}\tag{70}$$

which is null if \mathbf{B} is a solution of the equation:

$$\sum_{k=1}^K w_k \left(\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} = \left(\mathbf{B}^2 + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}}.\tag{71}$$

Therefore, for any probability measure $\mu \in \mathcal{G}$:

$$\begin{aligned}\left\langle \sum_{k=1}^K w_k \nabla_\beta S_\sigma(\alpha_k, \beta), \mu - \beta \right\rangle &= \left\langle \sum_{k=1}^K w_k g_k - h_\beta, \mu - \beta \right\rangle \\ &= \left\langle \sum_{k=1}^K w_k g_k(0) - h_\beta, \mu - \beta \right\rangle \\ &= \left\langle \sum_{w=1}^K w_k g_k(0) - h_\beta(0), \mu - \beta \right\rangle \\ &= \left(\sum_{w=1}^K w_k g_k(0) - h_\beta(0) \right) \int (\text{d}\mu - \text{d}\beta) \\ &= 0\end{aligned}\tag{72}$$

since both measures integrate to 1. Therefore, the optimality condition holds.

To end the proof, all we need to show is that (71) admits a positive definite solution. To show the existence of a solution, the same proof of Agueh and Carlier [1] applies. Indeed, let λ_k and Λ_k denote respectively the smallest and largest eigenvalue of \mathbf{A}_k . Let $\lambda = \min_k \lambda_k$ and $\Lambda = \max_k \Lambda_k$. Let $K_{\lambda, \Lambda}$ be the convex compact subset of positive definite matrices \mathbf{B} such that $\Lambda \text{Id} \succeq \mathbf{B} \succeq \lambda \text{Id}$. Define the map:

$$T : K_{\lambda, \Lambda} \rightarrow \mathcal{S}_{++}^d$$

$$\mathbf{B} \mapsto \left(\left(\sum_{k=1}^K w_k (\mathbf{B}^{\frac{1}{2}} \mathbf{A}_k \mathbf{B}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id})^{\frac{1}{2}} \right)^2 - \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}}$$

Now for any $\mathbf{B} \in K_{\lambda, \Lambda}$, it holds:

$$\lambda \text{Id} \preceq T(\mathbf{B}) \preceq \Lambda \text{Id}. \quad (73)$$

T is therefore a continuous function that maps $K_{\lambda, \Lambda}$ to itself, thus Brouwer's fixed-point theorem guarantees the existence of a solution. \square

Proof of Proposition 7

Proof. Using Fubini-Tonelli along with the optimality conditions (35), the double integral can be written:

$$\begin{aligned} \pi(\mathbb{R}^d \times \mathbb{R}^d) &= \int_{\mathbb{R}^d \times \mathbb{R}^d} e^{\frac{-\|x-y\|^2 + f(x) + g(y)}{2\sigma^2}} d\alpha(x) d\beta(y) \\ &= \int_{\mathbb{R}^d} \left(\int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2 + f(x)}{2\sigma^2}} d\alpha(x) \right) e^{\frac{g(y)}{2\sigma^2}} d\beta(y) \\ &= \int_{\mathbb{R}^d} e^{\frac{g(y)}{2\sigma^2} (1 - \frac{1}{\tau})} d\beta(y) \\ &= \int_{\mathbb{R}^d} e^{-\frac{g(y)}{\gamma}} d\beta(y) \end{aligned}$$

And similarly: $\pi(\mathbb{R}^d \times \mathbb{R}^d) = \int_{\mathbb{R}^d} e^{-\frac{f(x)}{\gamma}} d\alpha(x)$. Therefore, the three integrals in the dual objective (34) are equal to $\pi(\mathbb{R}^d \times \mathbb{R}^d)$ which ends the proof. \square

Lemma 3. [Sum of factorized quadratic forms] Let $\mathbf{A}, \mathbf{B} \in S_d^d$ such that $\mathbf{A} \neq \mathbf{B}$ and $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$. Denote $\alpha = (\mathbf{A}, \mathbf{a})$ and $\beta = (\mathbf{B}, \mathbf{b})$. Let $P_\alpha(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a})$ and $P_\beta(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{b})^\top \mathbf{B}(\mathbf{x} - \mathbf{b})$. Then:

$$P_\alpha(x) + P_\beta(x) = -\frac{1}{2} ((\mathbf{x} - \mathbf{c})^\top \mathbf{C}(\mathbf{x} - \mathbf{c}) + q_{\alpha, \beta}) \quad (74)$$

where:

$$\begin{cases} \mathbf{C} &= \mathbf{A} + \mathbf{B} \\ (\mathbf{A} + \mathbf{B})\mathbf{c} &= (\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b}) \\ q_{\alpha, \beta} &= \mathbf{a}^\top \mathbf{A}\mathbf{a} + \mathbf{b}^\top \mathbf{B}\mathbf{b} - \mathbf{c}^\top \mathbf{C}\mathbf{c} \end{cases} \quad (75)$$

In particular, if $\mathbf{C} = \mathbf{A} + \mathbf{B}$ is invertible, then:

$$\begin{cases} \mathbf{c} &= \mathbf{C}^{-1}(\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b}) \\ \mathbf{c}^\top \mathbf{C}\mathbf{c} &= (\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b})^\top \mathbf{C}^{-1}(\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b}) \end{cases} \quad (76)$$

Proof. On one hand,

$$\begin{aligned} P_\alpha(x) + P_\beta(x) &= -\frac{1}{2} ((\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a}) + (\mathbf{x} - \mathbf{b})^\top \mathbf{B}(\mathbf{x} - \mathbf{b})) \\ &= -\frac{1}{2} (\mathbf{x}^\top (\mathbf{A} + \mathbf{B})\mathbf{x} - 2\mathbf{x}^\top (\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b}) + \mathbf{a}^\top \mathbf{A}\mathbf{a} + \mathbf{b}^\top \mathbf{B}\mathbf{b}) \end{aligned}$$

On the other hand, for an arbitrary $\gamma = (\mathbf{c}, \mathbf{C})$ and $q \in \mathbb{R}$:

$$\begin{aligned} P_\gamma(x) - \frac{q}{2} &= -\frac{1}{2} ((\mathbf{x} - \mathbf{c})^\top \mathbf{C} (\mathbf{x} - \mathbf{C}) + q) \\ &= -\frac{1}{2} (x^\top \mathbf{C} x - 2x^\top \mathbf{C} \mathbf{c} + \mathbf{c}^\top \mathbf{C} \mathbf{c} + q) \end{aligned}$$

If $\mathbf{A} \neq \mathbf{B}$, identification of the parameters of both quadratic forms leads to (75). \square

Lemma 4. *[Gaussian convolution of factorized quadratic forms] Let $\mathbf{A} \in S_d$ and $\mathbf{a} \in \mathbb{R}^d$ and $\sigma > 0$ such that $\sigma^2 \mathbf{A} + \text{Id} \succ 0$. Let $Q_\alpha(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \mathbf{a})^\top \mathbf{A} (\mathbf{x} - \mathbf{a})$. Then the convolution of e^{Q_α} by the Gaussian kernel $\mathcal{N}(0, \frac{\text{Id}}{\sigma^2})$ is given by:*

$$\mathcal{N}(0, \frac{\text{Id}}{\sigma^2}) \star \exp(Q_\alpha) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \|\cdot - y\|^2 + Q_\alpha(y)\right) dy = c_\alpha \exp(Q(\mathbf{a}, \mathbf{J})) \quad (77)$$

where:

$$\begin{aligned} \mathbf{J} &= (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \\ c_\alpha &= \frac{1}{\sqrt{\det(\sigma^2 \mathbf{A} + \text{Id})}} \end{aligned}$$

Proof. Using Lemma 3 one can write for any $x \in \mathbb{R}^d$ considered fixed:

$$\begin{aligned} -\frac{1}{2\sigma^2} \|x - y\|^2 + Q_\alpha(y) &= Q(x, \frac{\text{Id}}{\sigma^2})(y) + Q(\mathbf{a}, \mathbf{A})(y) \\ &= Q(\mathbf{A}\mathbf{a} + \frac{x}{\sigma^2}, \mathbf{A} + \frac{\text{Id}}{\sigma^2})(y) + h(x) \end{aligned}$$

with $h(x) = -\frac{1}{2} (\mathbf{a}^\top \mathbf{A} \mathbf{a} + \frac{1}{\sigma^2} \|x\|^2 - \frac{1}{\sigma^2} (\sigma^2 \mathbf{A} \mathbf{a} + x)^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} (\sigma^2 \mathbf{A} \mathbf{a} + x))$. Therefore, the convolution integral is finite if and only if $\mathbf{A} + \frac{\text{Id}}{\sigma^2} \succ 0$ in which case we get the integral of a Gaussian density:

$$\begin{aligned} \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \int_{\mathbb{R}^d} \exp\left(Q(\mathbf{A}\mathbf{a} + \frac{x}{\sigma^2}, \mathbf{A} + \frac{\text{Id}}{\sigma^2})(y) + h(x)\right) dy &= \sqrt{\frac{\det(2\pi(\mathbf{A} + \frac{\text{Id}}{\sigma^2})^{-1})}{(2\pi\sigma^2)^n}} e^{h(x)} \\ &= \frac{e^{h(x)}}{\sqrt{\det(\sigma^2 \mathbf{A} + \text{Id})}} \end{aligned}$$

For the sake of clarity, let's separate the terms of h depending on their order in x : $h(x) = -\frac{1}{2} (h_2(x) + h_1(x) + h_0)$ where:

$$\begin{aligned} h_2(x) &= \frac{1}{\sigma^2} (\|x\|^2 - x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} x) \\ h_1(x) &= -2x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \mathbf{a} \\ h_0 &= \mathbf{a} \mathbf{A} \mathbf{a} - \sigma^2 \mathbf{a}^\top \mathbf{A} (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \mathbf{a} \end{aligned}$$

Finally, we can factorize h_2 and h_0 using Woodbury's matrix identity which holds even for a singular matrix \mathbf{A} :

$$(\sigma^2 \mathbf{A} + \text{Id})^{-1} = \text{Id} - \sigma^2 (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \quad (\text{Woodbury's identity})$$

Let $\mathbf{J} = (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A}$.

$$\begin{aligned}
h_2(x) &= \frac{1}{\sigma^2} (\|x\|^2 - x^\top (\text{Id} - \sigma^2 (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A}) x) \\
&= x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} x \\
&= x^\top \mathbf{J} x \\
h_1(x) &= -2x^\top \mathbf{J} \mathbf{a} \\
h_0 &= \mathbf{a} \mathbf{A} \mathbf{a} - \sigma^2 \mathbf{a}^\top \mathbf{A} (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \mathbf{a} \\
&= \mathbf{a}^\top \mathbf{A} (\text{Id} - \sigma^2 (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A}) \mathbf{a} \\
&= \mathbf{a}^\top \mathbf{A} (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{a} \\
&= \mathbf{a}^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \mathbf{a} \\
&= \mathbf{a}^\top \mathbf{J} \mathbf{a}
\end{aligned}$$

Therefore, $h(x) = -\frac{1}{2} (x^\top \mathbf{J} x - 2x^\top \mathbf{J} \mathbf{a} + \mathbf{a}^\top \mathbf{J} \mathbf{a}) = -\frac{1}{2} (x - \mathbf{a})^\top \mathbf{J} (x - \mathbf{a}) = \mathcal{Q}(\mathbf{a}, \mathbf{J})(x)$. \square

Lemma 5. [Gaussian convolution of generic quadratic forms] Let $\mathbf{A} \in S_d$ and $\mathbf{a} \in \mathbb{R}^d$ and $\sigma > 0$ such that $\sigma^2 \mathbf{A} + \text{Id} \succ 0$. Let $Q_\alpha(\mathbf{x}) = -\frac{1}{2} (\mathbf{x}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{a})$. Then the convolution of $e^{\mathcal{Q}_\alpha}$ by the Gaussian kernel $\mathcal{N}(0, \frac{\text{Id}}{\sigma^2})$ is given by:

$$\mathcal{N}(0, \frac{\text{Id}}{\sigma^2}) \star \exp(\mathcal{Q}_\alpha) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \|\cdot - y\|^2 + \mathcal{Q}_\alpha(y)\right) dy = c_\alpha \exp(\mathcal{Q}(\mathbf{G}\mathbf{a}, \mathbf{G}\mathbf{A})) \quad (78)$$

where:

$$\begin{aligned}
\mathbf{G} &= (\sigma^2 \mathbf{A} + \text{Id})^{-1} \\
c_\alpha &= \frac{e^{\frac{\sigma^2 \mathbf{a}^\top \mathbf{G} \mathbf{a}}{2}}}{\sqrt{\det(\sigma^2 \mathbf{A} + \text{Id})}}
\end{aligned}$$

Proof. Using Lemma 3 one can write for any $x \in \mathbb{R}^d$ considered fixed:

$$\begin{aligned}
-\frac{1}{2\sigma^2} \|x - y\|^2 + \mathcal{Q}_\alpha(y) &= \mathcal{Q}(x, \frac{\text{Id}}{\sigma^2})(y) + \mathcal{Q}(\mathbf{a}, \mathbf{A})(y) \\
&= \mathcal{Q}(\mathbf{a} + \frac{x}{\sigma^2}, \mathbf{A} + \frac{\text{Id}}{\sigma^2})(y) - \frac{1}{2\sigma^2} \|x\|^2 \\
&= \mathcal{Q}f((\sigma \mathbf{a} + \frac{x}{\sigma^2}, \mathbf{A} + \frac{\text{Id}}{\sigma^2})(y) + h(x)
\end{aligned}$$

with $h(x) = -\frac{1}{2} (\frac{1}{\sigma^2} \|x\|^2 - \frac{1}{\sigma^2} (\sigma^2 \mathbf{a} + x)^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} (\sigma^2 \mathbf{a} + x))$. Therefore, the convolution integral is finite if and only if $\mathbf{A} + \frac{\text{Id}}{\sigma^2} \succ 0$ in which case we get the integral of a Gaussian density:

$$\begin{aligned}
\frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \int_{\mathbb{R}^d} \exp\left(\mathcal{Q}f(\mathbf{a} + \frac{x}{\sigma^2}, \mathbf{A} + \frac{\text{Id}}{\sigma^2})(y) + h(x)\right) d(y) &= \sqrt{\frac{\det(2\pi(\mathbf{A} + \frac{\text{Id}}{\sigma^2})^{-1})}{(2\pi\sigma^2)^n}} e^{h(x)} \\
&= \frac{e^{h(x)}}{\sqrt{\det(\sigma^2 \mathbf{A} + \text{Id})}}
\end{aligned}$$

For the sake of clarity, let's separate the terms of h depending on their order in x : $h(x) = -\frac{1}{2} (h_2(x) + h_1(x) + h_0)$ where:

$$\begin{aligned}
h_2(x) &= \frac{1}{\sigma^2} (\|x\|^2 - x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} x) \\
h_1(x) &= -2x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{a} \\
h_0 &= -\sigma^2 \mathbf{a}^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{a}
\end{aligned}$$

Finally, we can factorize h_2 and h_0 using Woodbury's matrix identity which holds even for a singular matrix \mathbf{A} :

$$(\sigma^2 \mathbf{A} + \text{Id})^{-1} = \text{Id} - \sigma^2 (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} \quad (\text{Woodbury's identity})$$

Let $\mathbf{G} = (\sigma^2 \mathbf{A} + \text{Id})^{-1}$.

$$\begin{aligned}
h_2(x) &= \frac{1}{\sigma^2} (\|x\|^2 - x^\top (\text{Id} - \sigma^2 (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A}) x) \\
&= x^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{A} x \\
&= x^\top \mathbf{G} \mathbf{A} x \\
h_1(x) &= -2x^\top \mathbf{G} \mathbf{a} \\
h_0 &= -\sigma^2 \mathbf{a}^\top (\sigma^2 \mathbf{A} + \text{Id})^{-1} \mathbf{a} \\
&= -\sigma^2 \mathbf{a}^\top \mathbf{G} \mathbf{a}
\end{aligned}$$

Therefore, $h(x) = -\frac{1}{2} (x^\top \mathbf{G} \mathbf{A} x - 2x^\top \mathbf{G} \mathbf{a} - \sigma^2 \mathbf{a}^\top \mathbf{G} \mathbf{a}) = \mathcal{Q}(\mathbf{G} \mathbf{a}, \mathbf{G} \mathbf{A})(x) + \frac{\sigma^2 \mathbf{a}^\top \mathbf{G} \mathbf{a}}{2}$. \square

5.4 Proof of theorem 3

In the balanced case, we showed that Sinkhorn's transform is stable for quadratic potentials and that the resulting sequence is a contraction. Similarly, the following proposition shows that the unbalanced Sinkhorn transform is stable for quadratic potentials. \mathbf{M}

Proposition 8. *Let α be an unbalanced Gaussians given by $m_\alpha \mathcal{N}(\mathbf{a}, \mathbf{A})$. Let $\tau = \frac{\gamma}{2\sigma^2 + \gamma}$. Define the unbalanced Sinkhorn transform $T : \mathbb{R}^{\mathbb{R}^d} \rightarrow \mathbb{R}^{\mathbb{R}^d}$:*

$$T_\alpha(h)(x) \stackrel{\text{def}}{=} -\tau \log \int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2}{2\sigma^2} + h(y)} d\alpha(y) \quad (79)$$

Let $\mathbf{U} \in \mathcal{S}_d$, $\mathbf{u} \in \mathbb{R}^d$ and $m_u > 0$. If $h = \log(m_u) + \mathcal{Q}(\mathbf{u}, \mathbf{U})$ i.e $h(x) = \log(m_u) - \frac{1}{2}(x^\top \mathbf{U} x - 2x^\top \mathbf{u})$, then $T_\alpha(h)$ is well defined if and only if $\mathbf{F} \stackrel{\text{def}}{=} \sigma^2 \mathbf{U} + \sigma^2 \mathbf{A}^{-1} + \text{Id} \succ 0$, in which case $T_\alpha(h) = \mathcal{Q}(\mathbf{v}, \mathbf{V}) + \log(m_v)$ with the identified parameters:

$$\mathbf{V} = \tau \frac{1}{\sigma^2} (\mathbf{F}^{-1} - \text{Id}) \quad (80)$$

$$\mathbf{v} = -\tau \mathbf{F}^{-1} (\mathbf{A}^{-1} \mathbf{a} + \mathbf{u}) \quad (81)$$

$$m_v = \left(\frac{\sqrt{\det(\mathbf{A}) \det(\mathbf{F})}}{m_u m_\alpha e^{\frac{q_{u,\alpha}}{2}} \sigma^{2d}} \right)^\tau \quad (82)$$

where $q_{u,\alpha} = \frac{\sigma^2}{\tau^2} \mathbf{v}^\top \mathbf{F} \mathbf{v} - \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a}$.

Proof. The exponent inside the integral can be written as:

$$\begin{aligned}
e^{\frac{-\|x-y\|^2}{2\sigma^2} + h(y)} d\alpha(y) &\propto e^{\frac{-\|x-y\|^2}{2\sigma^2} - \frac{1}{2}(y^\top \mathbf{X} y - y^\top \mathbf{A}^{-1} y)} dy \\
&\propto e^{-\frac{1}{2}(y^\top (\frac{\text{Id}}{\sigma^2} + \mathbf{X} + \mathbf{A}^{-1}) y) + \frac{x^\top y}{\sigma^2}} dy
\end{aligned}$$

which is integrable if and only if $\mathbf{U} + \mathbf{A}^{-1} + \frac{1}{\sigma^2} \text{Id} \succ 0 \Leftrightarrow \mathbf{F} \succ 0$. Moreover, up to a multiplicative factor, the exponentiated Sinkhorn transform is equivalent to a Gaussian convolution of an

exponentiated quadratic form. Lemma 5 applies:

$$\begin{aligned}
e^{-T_\alpha(h)} &= \int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2}{2\sigma^2} + f(y)} d\alpha(y) \\
&= m_u m_\alpha \frac{\exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(2\pi \mathbf{A})}} \int_{\mathbb{R}^d} e^{\frac{-\|x-y\|^2}{2\sigma^2} + \mathcal{Q}(\mathbf{u}, \mathbf{U})(y) + \mathcal{Q}(\mathbf{A}^{-1} \mathbf{a}, \mathbf{A}^{-1})(y)} dy \\
&= m_u m_\alpha \frac{\exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(2\pi \mathbf{A})}} \sqrt{(2\pi\sigma^2)^{2d}} \exp(\mathcal{N}(\sigma^2 \text{Id})) \star \exp(\mathcal{Q}(\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}, \mathbf{U} + \mathbf{A}^{-1})) \\
&= m_u m_\alpha \frac{\sigma^{2d} \exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(\mathbf{A})}} \exp(\mathcal{N}(\sigma^2 \text{Id})) \star \exp(\mathcal{Q}(\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}, \mathbf{U} + \mathbf{A}^{-1})) \\
&= m_u m_\alpha \frac{\sigma^{2d} \exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(\mathbf{A})}} c_\alpha \exp(\mathcal{Q}(\mathbf{F}^{-1}(\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}), \mathbf{F}^{-1}(\mathbf{U} + \mathbf{A}^{-1}))) \\
&= m_u m_\alpha \frac{\sigma^{2d} \exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(\mathbf{A})}} c_\alpha \exp\left(\mathcal{Q}(\mathbf{F}^{-1}(\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}), \frac{1}{\sigma^2} \mathbf{F}^{-1}(\mathbf{F} - \text{Id}))\right) \\
&= m_u m_\alpha \frac{\sigma^{2d} \exp(-\frac{1}{2} \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a})}{\sqrt{\det(\mathbf{A})}} c_\alpha \exp\left(\mathcal{Q}(\mathbf{F}^{-1}(\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}), \frac{1}{\sigma^2} (\text{Id} - \mathbf{F}^{-1}))\right)
\end{aligned}$$

where $c_\alpha = \frac{\exp(\frac{1}{2} \sigma^2 (\mathbf{u} + \mathbf{A}^{-1} \mathbf{a})^\top \mathbf{F}^{-1} (\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}))}{\sqrt{\det(\mathbf{F})}}$.

Therefore, by applying $-\tau \log$ we can identify \mathbf{V} and \mathbf{v} . Substituting $\mathbf{u} + \mathbf{A}^{-1} \mathbf{a}$ by $-\frac{1}{\tau} \mathbf{F} \mathbf{v}$ leads to the equation of m_v . \square

Unlike the balanced case, the unbalanced Sinkhorn iterations require 2 more parameters (\mathbf{v} and m_v) with tangled updates. Proving the convergence of the resulting algorithm is more challenging. Instead, we directly solve the optimality conditions and show that a pair of quadratic potentials verifies (35).

Proposition 9. *The pair of quadratic forms (f, g) of (38) verifies the optimality conditions (35) if and only if:*

$$\begin{aligned}
\mathbf{F} &\stackrel{\text{def}}{=} \sigma^2 \mathbf{A}^{-1} + \sigma^2 \mathbf{U} + \text{Id} \succ 0 \\
\mathbf{G} &\stackrel{\text{def}}{=} \sigma^2 \mathbf{B}^{-1} + \sigma^2 \mathbf{V} + \text{Id} \succ 0,
\end{aligned} \tag{83}$$

$$\begin{aligned}
m_v \left(\frac{m_u m_\alpha e^{\frac{q_{u,\alpha}}{2}} \sigma^d}{\sqrt{\det(\mathbf{A}) \det(\mathbf{F})}} \right)^\tau &= 1 & m_u \left(\frac{m_v m_\beta e^{\frac{q_{v,\beta}}{2}} \sigma^d}{\sqrt{\det(\mathbf{B}) \det(\mathbf{G})}} \right)^\tau &= 1 \\
\mathbf{v} &= -\tau \mathbf{F}^{-1} (\mathbf{A}^{-1} \mathbf{a} + \mathbf{u}) & \mathbf{u} &= -\tau \mathbf{G}^{-1} (\mathbf{B}^{-1} \mathbf{b} + \mathbf{v}) \\
\mathbf{G} &= \tau \mathbf{F}^{-1} + \sigma^2 \mathbf{B}^{-1} + (1 - \tau) \text{Id} & \mathbf{F} &= \tau \mathbf{G}^{-1} + \sigma^2 \mathbf{A}^{-1} + (1 - \tau) \text{Id} \\
q_{u,\alpha} &= \frac{\sigma^2}{\tau^2} \mathbf{v}^\top \mathbf{F} \mathbf{v} - \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} & q_{v,\beta} &= \frac{\sigma^2}{\tau^2} \mathbf{u}^\top \mathbf{G} \mathbf{u} - \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b}
\end{aligned} \tag{84}$$

Proof. The equations on $m_u, m_v, \mathbf{u}, \mathbf{v}$ follow immediately from Proposition 8. Using the definition of \mathbf{F} and \mathbf{G} , substituting \mathbf{U} and \mathbf{V} leads to the equations in \mathbf{F} and \mathbf{G} . \square

We now turn to solve the system (84). Notice that in general, the dual potentials can only be identified up to an additive constant. Indeed, if a pair (f, g) is optimal, then $(f + K, g - K)$ is also optimal for any $K \in \mathbb{R}$ (the transportation plan does not change). Thus, at optimality, it is sufficient to obtain the product $m_u m_v$. We start by identifying (\mathbf{F}, \mathbf{G}) then (\mathbf{u}, \mathbf{v}) and finally $m_u m_v$.

Identifying \mathbf{F} and \mathbf{G} . The equations in \mathbf{F} and \mathbf{G} can shown to be equivalent to those of the balanced case up to some change of variables. Let $\lambda = \frac{1-\tau}{\sigma^2}$

$$\begin{aligned} & \begin{cases} \mathbf{F} &= \tau \mathbf{G}^{-1} + \sigma^2 \mathbf{A}^{-1} + (1-\tau) \text{Id} \\ \mathbf{G} &= \tau \mathbf{F}^{-1} + \sigma^2 \mathbf{B}^{-1} + (1-\tau) \text{Id} \end{cases} \\ & \Leftrightarrow \begin{cases} \mathbf{F} &= \left(\frac{\mathbf{G}}{\tau}\right)^{-1} + \frac{\sigma^2}{\tau} \tau (\mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id}) \\ \frac{\mathbf{G}}{\tau} &= \mathbf{F}^{-1} + \frac{\sigma^2}{\tau} (\mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id}) \end{cases} \\ & \Leftrightarrow \begin{cases} \mathbf{F} &= \tilde{\mathbf{G}}^{-1} + \sigma^2 \left(\frac{\tilde{\mathbf{A}}}{\tau}\right)^{-1} \\ \tilde{\mathbf{G}} &= \mathbf{F}^{-1} + \sigma^2 \tilde{\mathbf{B}}^{-1} \end{cases} \end{aligned}$$

which correspond to the balanced OT fixed point equations (20) associated with the pair $(\frac{\tilde{\mathbf{A}}}{\tau}, \tilde{\mathbf{B}})$ with the change of variables:

$$\tilde{\mathbf{G}} \stackrel{\text{def}}{=} \frac{\mathbf{G}}{\tau} \quad (85)$$

$$\tilde{\mathbf{A}} \stackrel{\text{def}}{=} \tau (\mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id})^{-1} \quad (86)$$

$$\tilde{\mathbf{B}} \stackrel{\text{def}}{=} \tau (\mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id})^{-1} \quad (87)$$

Notice that since $0 < \tau < 1$, $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{B}}$ are well-defined and positive definite. Therefore, Proposition 3 applies and we can write in closed form:

$$\begin{aligned} \mathbf{C} \stackrel{\text{def}}{=} \tilde{\mathbf{A}} \tilde{\mathbf{G}}^{-1} &= \left(\frac{1}{\tau} \tilde{\mathbf{A}} \tilde{\mathbf{B}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id} \\ &= \tilde{\mathbf{A}}^{\frac{1}{2}} \left(\frac{1}{\tau} \tilde{\mathbf{A}}^{\frac{1}{2}} \tilde{\mathbf{B}} \tilde{\mathbf{A}}^{\frac{1}{2}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \tilde{\mathbf{A}}^{-\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id} \end{aligned} \quad (88)$$

And similarly by symmetry:

$$\tilde{\mathbf{B}} \mathbf{F}^{-1} = \left(\frac{1}{\tau} \tilde{\mathbf{B}} \tilde{\mathbf{A}} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id} = \mathbf{C}^\top \quad (89)$$

Therefore we obtain \mathbf{F} and \mathbf{G} in closed form:

$$\mathbf{F} = \tilde{\mathbf{B}} \mathbf{C}^{-1} \quad (90)$$

$$\mathbf{G} = \mathbf{C}^{-1} \tilde{\mathbf{A}} \quad (91)$$

Finally, to obtain the formulas of $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{B}}$ of Theorem 3, use Woodburry's identity to write:

$$\begin{aligned} \tilde{\mathbf{B}} &= \tau \lambda (\text{Id} - \lambda (\mathbf{B} + \lambda \text{Id})^{-1}) \\ &= \frac{\gamma}{\gamma + 2\sigma^2} \frac{2\sigma^2 + \gamma}{2} (\text{Id} - \lambda (\mathbf{B} + \lambda \text{Id})^{-1}) \\ &= \frac{\gamma}{2} (\text{Id} - \lambda (\mathbf{B} + \lambda \text{Id})^{-1}) \end{aligned}$$

the same applies for $\tilde{\mathbf{A}}$.

Identifying \mathbf{u} and \mathbf{v} . Combining the equations in \mathbf{u} and \mathbf{v} leads to:

$$\begin{aligned} \mathbf{v} &= -\tau \mathbf{F}^{-1} (\mathbf{A}^{-1} \mathbf{a} + \tau \mathbf{u}) \\ \Leftrightarrow \mathbf{F} \mathbf{v} &= -\tau \mathbf{A}^{-1} \mathbf{a} - \tau \mathbf{u} \\ \Leftrightarrow \mathbf{F} \mathbf{v} &= -\tau \mathbf{A}^{-1} \mathbf{a} + \tau^2 \mathbf{G}^{-1} (\mathbf{B}^{-1} \mathbf{b} + \mathbf{v}) \\ \Leftrightarrow \mathbf{G} \mathbf{F} \mathbf{v} &= -\tau \mathbf{G} \mathbf{A}^{-1} \mathbf{a} + \tau^2 (\mathbf{B}^{-1} \mathbf{b} + \mathbf{v}) \\ \Leftrightarrow (\mathbf{G} \mathbf{F} - \tau^2 \text{Id}) \mathbf{v} &= -\tau \mathbf{G} \mathbf{A}^{-1} \mathbf{a} + \tau^2 \mathbf{B}^{-1} \mathbf{b} \end{aligned}$$

Similarly, $(\mathbf{F}\mathbf{G} - \tau^2 \text{Id})\mathbf{u} = -\tau\mathbf{F}\mathbf{B}^{-1}\mathbf{b} + \tau^2\mathbf{A}^{-1}\mathbf{a}$. Moreover, since $0 < \tau < 1$, it holds $(\mathbf{F} - \tau^2\mathbf{G}^{-1}) \succ (\mathbf{F} - \tau\mathbf{G}^{-1}) = \sigma^2\tilde{\mathbf{A}}^{-1} \succ 0$. Therefore, $(\mathbf{F}\mathbf{G} - \tau^2 \text{Id}) = (\mathbf{F} - \tau^2\mathbf{G}^{-1} \text{Id})\mathbf{G}$ is invertible. The same applies for $(\mathbf{G}\mathbf{F} - \tau^2 \text{Id})$.

Finally, both equations can be vectorized:

$$\begin{pmatrix} \mathbf{G}\mathbf{F} - \tau^2 \text{Id} & 0 \\ 0 & \mathbf{F}\mathbf{G} - \tau^2 \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{u} \end{pmatrix} = \begin{pmatrix} -\tau\mathbf{G} & \tau^2 \text{Id} \\ \tau^2 \text{Id} & -\tau\mathbf{F} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \quad (92)$$

Identifying $m_u m_v$. Now that \mathbf{F} , \mathbf{G} , \mathbf{u} and \mathbf{v} are given in closed form, $m_u m_v$ is obtained by taking the product of both equations:

$$(m_u m_v)^{\tau+1} = \left(\frac{\sqrt{\det(\mathbf{A}\mathbf{B}) \det(\mathbf{F}\mathbf{G})}}{\sigma^{2d} m_\alpha m_\beta} \right)^\tau \exp\left(-\frac{\tau}{2}(q_{u,\alpha} + q_{v,\beta})\right) \quad (93)$$

Transportation plan. Let $\omega \stackrel{\text{def}}{=} \frac{m_\alpha m_\beta}{\sqrt{\det(4\pi^2 \mathbf{A}\mathbf{B})}} m_u m_v e^{-\frac{1}{2}(\mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b})}$. At optimality, the transport plan π is given by:

$$\begin{aligned} \frac{d\pi}{dx dy}(x, y) &= \exp\left(\frac{f(x) + g(y) - \|x - y\|^2}{2\sigma^2}\right) \frac{d\alpha}{dx}(x) \frac{d\beta}{dy}(y) \\ &= \omega \exp\left(\mathcal{Q}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{u}, \mathbf{A}^{-1} + \mathbf{U})(x) - \frac{\|x - y\|^2}{2\sigma^2} + \mathcal{Q}(\mathbf{B}^{-1}\mathbf{b} + \mathbf{v}, \mathbf{B}^{-1} + \mathbf{V})(y)\right) \\ &= \omega \exp\left(\mathcal{Q}(\mathbf{U} + \mathbf{A}^{-1})(x) + \mathcal{Q}(\mathbf{V} + \mathbf{B}^{-1})(y) + \mathcal{Q}\left(-\frac{\text{Id}}{\sigma^2}, -\frac{\text{Id}}{\sigma^2}\right)(x, y)\right) \\ &= \omega \exp\left(\mathcal{Q}\left(\begin{pmatrix} \mathbf{A}^{-1}\mathbf{a} + \mathbf{u} \\ \mathbf{B}^{-1}\mathbf{b} + \mathbf{v} \end{pmatrix}, \begin{pmatrix} \mathbf{U} + \mathbf{A}^{-1} + \frac{\text{Id}}{\sigma^2} & 0 \\ 0 & \mathbf{V} + \mathbf{B}^{-1} + \frac{\text{Id}}{\sigma^2} \end{pmatrix}\right)(x, y)\right) \\ &= \omega \exp\left(\mathcal{Q}\left(\begin{pmatrix} \mathbf{A}^{-1}\mathbf{a} + \mathbf{u} \\ \mathbf{B}^{-1}\mathbf{b} + \mathbf{v} \end{pmatrix}, \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix}\right)(x, y)\right) \\ &= \omega \exp(\mathcal{Q}(\mu, \Gamma)(x, y)) \end{aligned}$$

with $\mu \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{A}^{-1}\mathbf{a} + \mathbf{u} \\ \mathbf{B}^{-1}\mathbf{b} + \mathbf{v} \end{pmatrix}$ and $\Gamma \stackrel{\text{def}}{=} \begin{pmatrix} \frac{\mathbf{F}}{\sigma^2} & -\frac{\text{Id}}{\sigma^2} \\ -\frac{\text{Id}}{\sigma^2} & \frac{\mathbf{G}}{\sigma^2} \end{pmatrix}$. Let's show that $\Gamma \succ 0$. Since $\frac{\mathbf{G}}{2\sigma^2} \succ 0$, it is sufficient to show that Schur complement $\frac{\mathbf{F}}{\sigma^2} - \frac{1}{\sigma^2}\mathbf{G}^{-1} \succ 0$. On one hand, with

$$\frac{\mathbf{F} - \mathbf{G}^{-1}}{\sigma^2} = \tau\tilde{\mathbf{A}}^{-1} - \frac{1}{\lambda}\mathbf{G}^{-1}$$

On the other hand, almost by definition $\tilde{\mathbf{A}} \prec \tau\lambda \text{Id}$ and $\tilde{\mathbf{B}} \prec \tau\lambda \text{Id}$. Thus for any $x \in \mathbb{R}^d$:

$$x^\top \frac{\tilde{\mathbf{A}}^{\frac{1}{2}} \tilde{\mathbf{B}} \tilde{\mathbf{A}}^{\frac{1}{2}}}{\tau} x \leq \lambda \|\tilde{\mathbf{A}}^{\frac{1}{2}} x\|^2 = \lambda x^\top \tilde{\mathbf{A}} x \leq \tau\lambda^2 \|x\|^2,$$

which implies

$$\left(\frac{\tilde{\mathbf{A}}^{\frac{1}{2}} \tilde{\mathbf{B}} \tilde{\mathbf{A}}^{\frac{1}{2}}}{\tau} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} \prec \sqrt{\tau\lambda^2 + \frac{\sigma^4}{4}} \text{Id} = \frac{\lambda}{2}(\sqrt{4\tau + (1-\tau)^2}) \text{Id} = \frac{\lambda(1+\tau)}{2} \text{Id}.$$

Therefore, using the second equality of (88) and inverting (90) to obtain \mathbf{G}^{-1} :

$$\begin{aligned}
x^\top \mathbf{G}^{-1} x &= x^\top \tilde{\mathbf{A}}^{-\frac{1}{2}} \left(\left(\frac{\tilde{\mathbf{A}}^{\frac{1}{2}} \tilde{\mathbf{B}} \tilde{\mathbf{A}}^{\frac{1}{2}}}{\tau} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \frac{\sigma^2}{2} \text{Id} \right) \tilde{\mathbf{A}}^{-\frac{1}{2}} x \\
&= (\tilde{\mathbf{A}}^{-\frac{1}{2}} x)^\top \left(\left(\frac{\tilde{\mathbf{A}}^{\frac{1}{2}} \tilde{\mathbf{B}} \tilde{\mathbf{A}}^{\frac{1}{2}}}{\tau} + \frac{\sigma^4}{4} \text{Id} \right)^{\frac{1}{2}} - \frac{\lambda(1-\tau)}{2} \text{Id} \right) (\tilde{\mathbf{A}}^{-\frac{1}{2}} x) \\
&\leq (\tilde{\mathbf{A}}^{-\frac{1}{2}} x)^\top \left(\frac{\lambda(1+\tau)}{2} \text{Id} - \frac{\lambda(1-\tau)}{2} \text{Id} \right) (\tilde{\mathbf{A}}^{-\frac{1}{2}} x) \\
&= \tau \lambda x^\top \tilde{\mathbf{A}}^{-1} x.
\end{aligned}$$

Thus $\mathbf{G}^{-1} \prec \tau \lambda \tilde{\mathbf{A}}^{-1}$. We can therefore conclude that the Schur complement $\frac{1}{\sigma^2}(\mathbf{F} - \mathbf{G}^{-1})$ is positive definite. By completing the square, we can factor $\frac{d\pi}{dx dy}$ as a Gaussian density. Let $z \stackrel{\text{def}}{=} \begin{pmatrix} x \\ y \end{pmatrix}$:

$$\begin{aligned}
\frac{d\pi}{dx dy}(x, y) &= \omega \exp(\mathcal{Q}(\mu, \Gamma)(x, y)) \\
&= \omega \exp\left(-\frac{1}{2}(z^\top \Gamma z - 2z^\top \mu)\right) \\
&= \omega \exp\left(\frac{1}{2}\mu^\top \Gamma^{-1} \mu - \frac{1}{2}(z - \Gamma^{-1} \mu)^\top \Gamma (z - \Gamma^{-1} \mu)\right) \\
&= \omega e^{\frac{1}{2}\mu^\top \Gamma^{-1} \mu} \mathcal{N}(\mathbf{H}\mu, \mathbf{H})(z),
\end{aligned}$$

where $\mathbf{H} = \Gamma^{-1}$.

Detailed expressions. To conclude the proof of Theorem 3, we need to simplify the formulas of m , $\mathbf{H}\mu$ and \mathbf{H} . First, we will start with the mean $\mathbf{H}\mu$.

$\mathbf{H}\mu$ Using the optimality conditions of Proposition 9 and the closed form formula of \mathbf{v} and \mathbf{u} :

$$\begin{aligned}
\mu &= \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} + \mathbf{u} \\ \mathbf{B}^{-1} \mathbf{b} + \mathbf{v} \end{pmatrix} \\
&= -\frac{1}{\tau} \begin{pmatrix} \mathbf{F} \mathbf{v} \\ \mathbf{G} \mathbf{u} \end{pmatrix} \\
&= -\frac{1}{\tau} \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{u} \end{pmatrix} \\
&= -\frac{1}{\tau} \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix} \begin{pmatrix} \mathbf{G} \mathbf{F} - \tau^2 \text{Id} & 0 \\ 0 & \mathbf{F} \mathbf{G} - \tau^2 \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} -\tau \mathbf{G} & \tau^2 \text{Id} \\ \tau^2 \text{Id} & -\tau \mathbf{F} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix} \begin{pmatrix} \mathbf{G} \mathbf{F} - \tau^2 \text{Id} & 0 \\ 0 & \mathbf{F} \mathbf{G} - \tau^2 \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{G} & -\tau \text{Id} \\ -\tau \text{Id} & \mathbf{F} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \quad (94) \\
&= \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix} \begin{pmatrix} (\mathbf{F} - \tau^2 \mathbf{G}^{-1})^{-1} & -\tau(\mathbf{G} \mathbf{F} - \tau^2 \text{Id})^{-1} \\ -\tau(\mathbf{F} \mathbf{G} - \tau^2 \text{Id})^{-1} & (\mathbf{G} - \tau^2 \mathbf{F}^{-1})^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix} \begin{pmatrix} \mathbf{F} & \tau \text{Id} \\ \tau \text{Id} & \mathbf{G} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \begin{pmatrix} \text{Id} & \tau \mathbf{G}^{-1} \\ \tau \mathbf{F}^{-1} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}
\end{aligned}$$

Therefore:

$$\begin{aligned}
\mathbf{H}\mu &= \sigma^2 \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix}^{-1} \begin{pmatrix} \text{Id} & \tau\mathbf{G}^{-1}\text{Id} \\ \tau\mathbf{F}^{-1}\text{Id} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \sigma^2 \left(\begin{pmatrix} \text{Id} & \tau\mathbf{G}^{-1}\text{Id} \\ \tau\mathbf{F}^{-1}\text{Id} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix} \right)^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \sigma^2 \begin{pmatrix} \mathbf{F} - \tau\mathbf{G}^{-1} & -(1-\tau)\text{Id} \\ -(1-\tau)\text{Id} & \mathbf{G} - \tau\mathbf{F}^{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \sigma^2 \begin{pmatrix} \sigma^2\mathbf{A}^{-1} + (1-\tau)\text{Id} & -(1-\tau)\text{Id} \\ -(1-\tau)\text{Id} & \sigma^2\mathbf{B}^{-1} + (1-\tau)\text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{A}^{-1} + \text{Id} & -\lambda\text{Id} \\ -\lambda\text{Id} & \mathbf{B}^{-1} + \lambda\text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}
\end{aligned} \tag{95}$$

Let's compute the inverse of:

$$\mathbf{Z} \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{A}^{-1} + \frac{1}{\lambda}\text{Id} & -\frac{1}{\lambda}\text{Id} \\ -\frac{1}{\lambda}\text{Id} & \mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id} \end{pmatrix}. \tag{96}$$

Let \mathbf{S} and \mathbf{S}' be the respective Schur complements of $\mathbf{A}^{-1} + \frac{1}{\lambda}\text{Id}$ and $\mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id}$ in \mathbf{Z} . The block inverse formula writes:

$$\mathbf{Z}^{-1} = \begin{pmatrix} \mathbf{S} & \frac{1}{\lambda}\mathbf{S}(\mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id})^{-1} \\ \frac{1}{\lambda}(\mathbf{A}^{-1} + \frac{1}{\lambda}\text{Id})^{-1}\mathbf{S} & \mathbf{S}' \end{pmatrix}.$$

Using Woodbury's identity twice and denoting $\mathbf{X} \stackrel{\text{def}}{=} \mathbf{A} + \mathbf{B} + \lambda\text{Id}$:

$$\begin{aligned}
\mathbf{S} &= (\mathbf{A}^{-1} + \frac{1}{\lambda}\text{Id} - \frac{1}{\lambda^2}(\mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id})^{-1})^{-1} \\
&= (\mathbf{A}^{-1} + (\mathbf{B} + \lambda\text{Id})^{-1})^{-1} \\
&= (\mathbf{A} - \mathbf{A}(\mathbf{A} + \mathbf{B} + \lambda\text{Id})^{-1}\mathbf{A}) \\
&= \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}.
\end{aligned}$$

And similarly: $\mathbf{S}' = \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B}$. The off-diagonal blocks can be simplified as well:

$$\begin{aligned}
\frac{1}{\lambda}\mathbf{S}(\mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id})^{-1} &= \frac{1}{\lambda}(\mathbf{A}^{-1} + (\mathbf{B} + \lambda\text{Id})^{-1})^{-1}(\mathbf{B}^{-1} + \frac{1}{\lambda}\text{Id})^{-1} \\
&= (\mathbf{A}^{-1} + (\mathbf{B} + \lambda\text{Id})^{-1})^{-1}(\lambda\text{Id} + \mathbf{B}\text{Id})^{-1}\mathbf{B} \\
&= ((\mathbf{B} + \lambda\text{Id}) - (\mathbf{B} + \lambda\text{Id})(\mathbf{A} + \mathbf{B} + \lambda\text{Id})^{-1}(\mathbf{B} + \lambda\text{Id}))(\lambda\text{Id} + \mathbf{B}\text{Id})^{-1}\mathbf{B} \\
&= \mathbf{B} - (\mathbf{B} + \lambda\text{Id})\mathbf{X}^{-1}\mathbf{B} \\
&= \mathbf{B} - (\mathbf{X} - \mathbf{A})\mathbf{X}^{-1}\mathbf{B} \\
&= \mathbf{A}\mathbf{X}^{-1}\mathbf{B}.
\end{aligned}$$

Similarly, $\frac{1}{\lambda}(\mathbf{A}^{-1} + \frac{1}{\lambda}\text{Id})^{-1}\mathbf{S} = \mathbf{B}\mathbf{X}^{-1}\mathbf{A}$. Thus, the inverse of \mathbf{Z} is given by:

$$\mathbf{Z}^{-1} = \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix}. \tag{97}$$

and finally:

$$\begin{aligned}
\mathbf{H}\mu &= \mathbf{Z}^{-1} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} = \begin{pmatrix} \text{Id} - \mathbf{A}\mathbf{X}^{-1} & \mathbf{A}\mathbf{X}^{-1} \\ \mathbf{B}\mathbf{X}^{-1} & \text{Id} - \mathbf{B}\mathbf{X}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{a} + \mathbf{A}\mathbf{X}^{-1}(\mathbf{b} - \mathbf{a}) \\ \mathbf{b} + \mathbf{B}\mathbf{X}^{-1}(\mathbf{a} - \mathbf{b}) \end{pmatrix}
\end{aligned}$$

Finding the covariance matrix \mathbf{H} . To compute $\mathbf{H} = \left(\frac{1}{\sigma^2} \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix} \right)^{-1}$ one may use the block inverse formula. However, the Schur complement $(\mathbf{F} - \mathbf{G}^{-1})^{-1}$ is not easy to manipulate. Instead notice that the following holds:

$$\begin{aligned} \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix} \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix} &= \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{F} - \tau \mathbf{G}^{-1} & -(1-\tau) \text{Id} \\ -(1-\tau) \text{Id} & \mathbf{G} - \tau \mathbf{F}^{-1} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id} & -\frac{1}{\lambda} \text{Id} \\ -\frac{1}{\lambda} \text{Id} & \mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id} \end{pmatrix}, \end{aligned}$$

where the last equality follows from the optimality conditions (84). Therefore:

$$\mathbf{H} = \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id} & -\frac{1}{\lambda} \text{Id} \\ -\frac{1}{\lambda} \text{Id} & \mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id} \end{pmatrix}^{-1}.$$

Notice that we have already computed the inverse matrix on the right side above in the developments of $\mathbf{H}\mu$. Thus:

$$\begin{aligned} \mathbf{H} &= \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \text{Id} & \tau \mathbf{C}\tilde{\mathbf{B}}^{-1} \\ \mathbf{C}^\top \tilde{\mathbf{A}}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \text{Id} & \mathbf{C}(\mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id}) \\ \mathbf{C}^\top(\mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id}) & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \text{Id} & \mathbf{C}(\mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id}) \\ \mathbf{C}^\top(\mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id}) & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \text{Id} & \frac{1}{\lambda} \mathbf{C}(\lambda \text{Id} + \mathbf{B})\mathbf{B}^{-1} \\ \frac{1}{\lambda} \mathbf{C}^\top \mathbf{C}(\lambda \text{Id} + \mathbf{A})\mathbf{A}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \text{Id} & \frac{1}{\lambda} \mathbf{C}(\mathbf{X} - \mathbf{A})\mathbf{B}^{-1} \\ \frac{1}{\lambda} \mathbf{C}^\top(\mathbf{X} - \mathbf{B})\mathbf{A}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A} + \frac{1}{\lambda} \mathbf{C}(\mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}) & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} + \frac{1}{\lambda} \mathbf{C}(\mathbf{X} - \mathbf{A})(\text{Id} - \mathbf{X}^{-1}\mathbf{B}) \\ \frac{1}{\lambda} \mathbf{C}^\top(\mathbf{X} - \mathbf{B})(\text{Id} - \mathbf{X}^{-1}\mathbf{A}) + \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \frac{1}{\lambda} \mathbf{C}^\top(\mathbf{X} - \mathbf{B})\mathbf{X}^{-1}\mathbf{B} + \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} (\text{Id} + \frac{1}{\lambda} \mathbf{C})(\mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}) & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} + \frac{1}{\lambda} \mathbf{C}(\mathbf{X} - \mathbf{A} - \mathbf{B} + \mathbf{A}\mathbf{X}^{-1}\mathbf{B}) \\ \lambda \mathbf{C}^\top(\lambda \text{Id} + \mathbf{B}\mathbf{X}^{-1}\mathbf{A}) + \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & \frac{1}{\lambda} \mathbf{C}^\top(\mathbf{X} - \mathbf{B})\mathbf{X}^{-1}\mathbf{B} + \mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \\ &= \begin{pmatrix} (\text{Id} + \frac{1}{\lambda} \mathbf{C})(\mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}) & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} + \frac{1}{\lambda} \mathbf{C}(\lambda \text{Id} + \mathbf{A}\mathbf{X}^{-1}\mathbf{B}) \\ \mathbf{C}^\top + \frac{1}{\lambda} \mathbf{C}^\top \mathbf{B}\mathbf{X}^{-1}\mathbf{A} + \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & (\text{Id} + \frac{1}{\lambda} \mathbf{C}^\top)(\mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B}) \end{pmatrix} \\ &= \begin{pmatrix} (\text{Id} + \frac{1}{\lambda} \mathbf{C})(\mathbf{A} - \mathbf{A}\mathbf{X}^{-1}\mathbf{A}) & \mathbf{C} + (\text{Id} + \frac{1}{\lambda} \mathbf{C})\mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{C}^\top + (\text{Id} + \frac{1}{\lambda} \mathbf{C}^\top)\mathbf{B}\mathbf{X}^{-1}\mathbf{A} & (\text{Id} + \frac{1}{\lambda} \mathbf{C}^\top)(\mathbf{B} - \mathbf{B}\mathbf{X}^{-1}\mathbf{B}) \end{pmatrix}. \end{aligned}$$

Finding the mass of the plan π . The optimal transport plan is given by:

$$\frac{d\pi}{dx dy}(x, y) = \omega e^{\frac{1}{2}\mu^\top \Gamma^{-1} \mu} \sqrt{\det(2\pi \mathbf{H})} \mathcal{N}(\mathbf{H}\mu, \mathbf{H})(z), \quad (98)$$

where

$$\begin{aligned} \omega &= \frac{m_\alpha m_\beta}{\sqrt{\det(4\pi^2 \mathbf{A}\mathbf{B})}} m_u m_v e^{-\frac{1}{2}(\mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b})} \\ &= \frac{m_\alpha m_\beta}{\sqrt{\det(4\pi^2 \mathbf{A}\mathbf{B})}} \left(\frac{\sqrt{\det(\mathbf{A}\mathbf{B}) \det(\mathbf{F}\mathbf{G})}}{\sigma^{2d} m_\alpha m_\beta} \right)^{\frac{\tau}{\tau+1}} e^{-\frac{\tau}{2(\tau+1)}(q_{u,\alpha} + q_{v,\beta})} e^{-\frac{1}{2}(\mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b})} \\ &= \frac{1}{(2\pi)^d} \left(\frac{m_\alpha m_\beta}{\sqrt{\det(\mathbf{A}\mathbf{B})}} \right)^{\frac{1}{\tau+1}} \left(\frac{\sqrt{\det(\mathbf{F}\mathbf{G})}}{\sigma^{2d}} \right)^{\frac{\tau}{\tau+1}} e^{-\frac{\tau}{2(\tau+1)}(q_{u,\alpha} + q_{v,\beta})} e^{-\frac{1}{2}(\mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b})}. \end{aligned}$$

First, let's simplify the argument of the exponential terms. Isolating the terms that depend only on the input means \mathbf{a}, \mathbf{b} it holds: $q_{u,\alpha} + q_{v,\beta} = \frac{\sigma^2}{\tau^2}(\mathbf{v}^\top \mathbf{F} \mathbf{v} + \mathbf{u}^\top \mathbf{G} \mathbf{u}) + \mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b}$. Therefore, the full exponential argument is given by:

$$\phi \stackrel{\text{def}}{=} \mu^\top \Gamma^{-1} \mu - \frac{\tau}{\tau+1} \frac{\sigma^2}{\tau^2} (\mathbf{v}^\top \mathbf{F} \mathbf{v} + \mathbf{u}^\top \mathbf{G} \mathbf{u}) - \frac{1}{\tau+1} (\mathbf{a}^\top \mathbf{A}^{-1} \mathbf{a} + \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b}) \quad (99)$$

On one hand, using Equation (95) we replace μ :

$$\begin{aligned} \mu^\top \Gamma^{-1} \mu &= \mu^\top \mathbf{H} \mu \\ &= \sigma^2 \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{F} & -\text{Id} \\ -\text{Id} & \mathbf{G} \end{pmatrix}^{-1} \begin{pmatrix} \text{Id} & \tau \mathbf{G}^{-1} \\ \tau \mathbf{F}^{-1} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix} \end{aligned}$$

On the other hand:

$$\begin{aligned} \frac{\sigma^2}{\tau^2} (\mathbf{v}^\top \mathbf{F} \mathbf{v} + \mathbf{u}^\top \mathbf{G} \mathbf{u}) &= \sigma^2 ((\mathbf{A}^{-1} \mathbf{a} + \mathbf{u})^\top \mathbf{F}^{-1} (\mathbf{A}^{-1} \mathbf{a} + \mathbf{u}) + (\mathbf{B}^{-1} \mathbf{b} + \mathbf{v})^\top \mathbf{G}^{-1} (\mathbf{B}^{-1} \mathbf{b} + \mathbf{v})) \\ &= \sigma^2 \mu^\top \begin{pmatrix} \mathbf{F}^{-1} & 0 \\ 0 & \mathbf{G}^{-1} \end{pmatrix} \mu \\ &= \sigma^2 \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{F}^{-1} & 0 \\ 0 & \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \text{Id} & \tau \mathbf{G}^{-1} \\ \tau \mathbf{F}^{-1} & \text{Id} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix} \end{aligned}$$

Let $\mathbf{J} = \begin{pmatrix} \text{Id} & \tau \mathbf{G}^{-1} \\ \tau \mathbf{F}^{-1} & \text{Id} \end{pmatrix}$ and $\mathbf{K} = \begin{pmatrix} \mathbf{F} & 0 \\ 0 & \mathbf{G} \end{pmatrix}$. It holds:

$$\mu^\top \Gamma^{-1} \mu - \frac{\tau}{\tau+1} \frac{\sigma^2}{\tau^2} (\mathbf{v}^\top \mathbf{F} \mathbf{v} + \mathbf{u}^\top \mathbf{G} \mathbf{u}) = \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}^\top \mathbf{J}^{\top-1} (\mathbf{H} - \frac{\sigma^2 \tau}{\tau+1} \mathbf{K}^{-1}) \mathbf{J}^{-1} \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}$$

Let's compute the matrix $\mathbf{J}^{\top-1} (\mathbf{H} - \frac{\sigma^2 \tau}{\tau+1} \mathbf{K}^{-1}) \mathbf{J}^{-1}$. First keep in mind that $\mathbf{J} \mathbf{K} = \begin{pmatrix} \mathbf{F} & \tau \text{Id} \\ \tau \text{Id} & \mathbf{G} \end{pmatrix}$.

Now using Woodbury's identity:

$$\begin{aligned} \left(\mathbf{J}^{\top-1} (\mathbf{H} - \frac{\tau}{\tau+1} \mathbf{K}^{-1}) \mathbf{J}^{-1} \right)^{-1} &= \mathbf{J} (\mathbf{H} - \frac{\tau \sigma^2}{\tau+1} \mathbf{K}^{-1})^{-1} \mathbf{J}^\top \\ &= \mathbf{J} \left(-\frac{\tau+1}{\tau \sigma^2} \mathbf{K} - \left(\frac{\tau+1}{\tau \sigma^2} \right)^2 \mathbf{K} (\mathbf{H}^{-1} - \frac{\tau+1}{\tau \sigma^2} \mathbf{K})^{-1} \mathbf{K} \right) \mathbf{J}^\top \\ &= \frac{\tau+1}{\tau \sigma^2} \left(-\mathbf{J} \mathbf{K} \mathbf{J}^\top - \frac{\tau+1}{\tau \sigma^2} \mathbf{J} \mathbf{K} \left(\begin{pmatrix} -\frac{\mathbf{F}}{\tau \sigma^2} & -\frac{1}{\sigma^2} \text{Id} \\ -\frac{1}{\sigma^2} \text{Id} & -\frac{\mathbf{G}}{\tau \sigma^2} \end{pmatrix}^{-1} (\mathbf{J} \mathbf{K}^\top)^\top \right) \right) \\ &= \frac{\tau+1}{\tau \sigma^2} \left(-\mathbf{J} \mathbf{K} \mathbf{J}^\top + (\tau+1) \mathbf{J} \mathbf{K} \left(\begin{pmatrix} \mathbf{F} & \tau \text{Id} \\ \tau \text{Id} & \mathbf{G} \end{pmatrix}^{-1} (\mathbf{J} \mathbf{K}^\top)^\top \right) \right) \\ &= \frac{\tau+1}{\tau \sigma^2} \left(-\begin{pmatrix} \mathbf{F} & \tau \text{Id} \\ \tau \text{Id} & \mathbf{G} \end{pmatrix} \begin{pmatrix} \text{Id} & \tau \mathbf{F}^{-1} \\ \tau \mathbf{G}^{-1} & \text{Id} \end{pmatrix} + (\tau+1) \begin{pmatrix} \mathbf{F} & \tau \text{Id} \\ \tau \text{Id} & \mathbf{G} \end{pmatrix} \right) \\ &= \frac{\tau+1}{\tau \sigma^2} \begin{pmatrix} -\mathbf{F} - \tau^2 \mathbf{G}^{-1} + (\tau+1) \mathbf{F} & (-2\tau + \tau(\tau+1)) \text{Id} \\ (-2\tau + \tau(\tau+1)) \text{Id} & -\mathbf{G} - \tau^2 \mathbf{F}^{-1} + (\tau+1) \mathbf{G} \end{pmatrix} \\ &= \frac{\tau+1}{\sigma^2} \begin{pmatrix} \mathbf{F} - \tau \mathbf{G}^{-1} & -(1-\tau) \text{Id} \\ -(1-\tau) \text{Id} & \mathbf{G} - \tau \mathbf{F}^{-1} \end{pmatrix} \\ &= (\tau+1) \begin{pmatrix} \mathbf{A}^{-1} + \frac{1}{\lambda} \text{Id} & -\frac{1}{\lambda} \text{Id} \\ -\frac{1}{\lambda} \text{Id} & \mathbf{B}^{-1} + \frac{1}{\lambda} \text{Id} \end{pmatrix} \\ &= (\tau+1) \mathbf{Z} \end{aligned}$$

Therefore:

$$\mu^\top \Gamma^{-1} \mu - \frac{\tau}{\tau+1} \frac{\sigma^2}{\tau^2} (\mathbf{v}^\top \mathbf{F} \mathbf{v} + \mathbf{u}^\top \mathbf{G} \mathbf{u}) = \frac{1}{\tau+1} \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}^\top \mathbf{Z}^{-1} \begin{pmatrix} \mathbf{A}^{-1} \mathbf{a} \\ \mathbf{B}^{-1} \mathbf{b} \end{pmatrix} \quad (100)$$

The full exponential argument ϕ defined in Equation (99) is given by:

$$\begin{aligned}
\phi &= \frac{1}{\tau+1} \left(\begin{pmatrix} \mathbf{A}^{-1}\mathbf{a} \\ \mathbf{B}^{-1}\mathbf{b} \end{pmatrix}^\top \mathbf{Z}^{-1} \begin{pmatrix} \mathbf{A}^{-1}\mathbf{a} \\ \mathbf{B}^{-1}\mathbf{b} \end{pmatrix} - \mathbf{a}^\top \mathbf{A}^{-1}\mathbf{a} - \mathbf{b}^\top \mathbf{B}^{-1}\mathbf{b} \right) \\
&= \frac{1}{\tau+1} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \left(\mathbf{Z}^{-1} - \begin{pmatrix} \mathbf{A} & 0 \\ 0 & \mathbf{B} \end{pmatrix} \right) \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \frac{1}{\tau+1} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} -\mathbf{A}\mathbf{X}^{-1}\mathbf{A} & \mathbf{A}\mathbf{X}^{-1}\mathbf{B} \\ \mathbf{B}\mathbf{X}^{-1}\mathbf{A} & -\mathbf{B}\mathbf{X}^{-1}\mathbf{B} \end{pmatrix} \begin{pmatrix} \mathbf{A}^{-1} & 0 \\ 0 & \mathbf{B}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= \frac{1}{\tau+1} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} -\mathbf{X}^{-1} & \mathbf{X}^{-1} \\ \mathbf{X}^{-1} & -\mathbf{X}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \\
&= -\frac{1}{\tau+1} (\mathbf{a} - \mathbf{b})^\top \mathbf{X}^{-1} (\mathbf{a} - \mathbf{b}) \\
&= \frac{1}{\tau+1} \|\mathbf{a} - \mathbf{b}\|_{\mathbf{X}^{-1}}^2
\end{aligned}$$

Substituting in (98) leads to:

$$\begin{aligned}
m_\pi &\stackrel{\text{def}}{=} \pi(\mathbb{R}^d \times \mathbb{R}^d) \\
&= \sqrt{\det(\mathbf{H})} \left(\frac{m_\alpha m_\beta}{\sqrt{\det(\mathbf{AB})}} \right)^{\frac{1}{\tau+1}} \left(\frac{\sqrt{\det(\mathbf{FG})}}{\sigma^{2d}} \right)^{\frac{\tau}{\tau+1}} e^{-\frac{1}{2(\tau+1)} (\|\mathbf{a} - \mathbf{b}\|_{\mathbf{X}^{-1}}^2)}.
\end{aligned}$$

The determinants can be easily expressed as functions of \mathbf{C} . First notice that:

$$\det(\mathbf{H}) = \frac{1}{\det(\Gamma)} = \frac{\sigma^{4d}}{\det(\mathbf{FG} - \text{Id})},$$

and using the definition of \mathbf{C} , it holds that

$$\mathbf{FG} = \tilde{\mathbf{B}}\mathbf{C}^{-2}\tilde{\mathbf{A}}.$$

Therefore, $\det(\mathbf{FG}) = \frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})}{\det(\mathbf{C})^2}$. Keeping in mind that the closed form expression of \mathbf{C} given in (90) is applied to the pair $(\frac{1}{\tau}\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ in the unbalanced case, it holds: $\mathbf{C}^2 + \sigma^2\mathbf{C} = \frac{1}{\tau}\tilde{\mathbf{A}}\tilde{\mathbf{B}}$. Thus:

$$\begin{aligned}
\mathbf{FG} - \text{Id} &= \tilde{\mathbf{B}}\mathbf{C}^{-2}\tilde{\mathbf{A}}(\text{Id} - \tilde{\mathbf{A}}^{-1}\mathbf{C}^2\tilde{\mathbf{B}}^{-1}) \\
&= \tilde{\mathbf{B}}\mathbf{C}^{-2}\tilde{\mathbf{A}}(\text{Id} - \tilde{\mathbf{A}}^{-1}(\frac{1}{\tau}\tilde{\mathbf{A}}\tilde{\mathbf{B}} - \sigma^2\mathbf{C})\tilde{\mathbf{B}}^{-1}) \\
&= \tilde{\mathbf{B}}\mathbf{C}^{-2}\tilde{\mathbf{A}}(\frac{(1-\tau)}{\tau}\text{Id} + \sigma^2\tilde{\mathbf{A}}^{-1}\mathbf{C}\tilde{\mathbf{B}}^{-1}) \\
&= \sigma^2\tilde{\mathbf{B}}\mathbf{C}^{-2}\tilde{\mathbf{A}}(-\frac{2}{\gamma}\text{Id} + \tilde{\mathbf{A}}^{-1}\mathbf{C}\tilde{\mathbf{B}}^{-1}) \\
&= \sigma^2\tilde{\mathbf{B}}\mathbf{C}^{-2}(-\frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}} + \mathbf{C})\tilde{\mathbf{B}}^{-1},
\end{aligned}$$

therefore

$$\det(\mathbf{FG} - \text{Id}) = \sigma^{2d} \frac{\det((-\frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}} + \mathbf{C})}{\det(\mathbf{C})^2}.$$

Replacing the determinant formulas of \mathbf{FG} and $\mathbf{FG} - \text{Id}$ and re-arranging the common terms $\det(\mathbf{C})$ and σ leads to:

$$\begin{aligned}
\pi(\mathbb{R}^d \times \mathbb{R}^d) &= \frac{\left(m_\alpha m_\beta \sigma^{2d} \det(\mathbf{C}) \sqrt{\frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})^\tau}{\det(\mathbf{A}\mathbf{B})}}\right)^{\frac{1}{\tau+1}}}{\sqrt{\frac{\det(\mathbf{C} - \frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}})}{\sigma^{2d}}}} e^{-\frac{1}{2(\tau+1)}(\|\mathbf{a}-\mathbf{b}\|_{\mathbf{X}-1}^2)} \\
&= \sigma^{d(\frac{2}{\tau+1}-1)} \frac{\left(m_\alpha m_\beta \det(\mathbf{C}) \sqrt{\frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})^\tau}{\det(\mathbf{A}\mathbf{B})}}\right)^{\frac{1}{\tau+1}}}{\sqrt{\det(\mathbf{C} - \frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}})}} e^{-\frac{1}{2(\tau+1)}(\|\mathbf{a}-\mathbf{b}\|_{\mathbf{X}-1}^2)} \\
&= \sigma^{d\frac{1-\tau}{\tau+1}} \frac{\left(m_\alpha m_\beta \det(\mathbf{C}) \sqrt{\frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})^\tau}{\det(\mathbf{A}\mathbf{B})}}\right)^{\frac{1}{\tau+1}}}{\sqrt{\det(\mathbf{C} - \frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}})}} e^{-\frac{1}{2(\tau+1)}(\|\mathbf{a}-\mathbf{b}\|_{\mathbf{X}-1}^2)} \\
&= \sigma^{\frac{d\sigma^2}{\sigma^2+\gamma}} \frac{\left(m_\alpha m_\beta \det(\mathbf{C}) \sqrt{\frac{\det(\tilde{\mathbf{A}}\tilde{\mathbf{B}})^\tau}{\det(\mathbf{A}\mathbf{B})}}\right)^{\frac{1}{\tau+1}}}{\sqrt{\det(\mathbf{C} - \frac{2}{\gamma}\tilde{\mathbf{A}}\tilde{\mathbf{B}})}} e^{-\frac{1}{2(\tau+1)}(\|\mathbf{a}-\mathbf{b}\|_{\mathbf{X}-1}^2)}
\end{aligned} \tag{101}$$

Deriving a closed form for UOT_σ . Using Equation (101), a direct application of Proposition 7 yields

$$\text{UOT}_\sigma(\alpha, \beta) = \gamma(m_\alpha + m_\beta) + 2\sigma^2(m_\alpha m_\beta) - 2(\sigma^2 + 2\gamma)m_{\pi^*}. \tag{102}$$

This ends the proof of Theorem 3.